

# Performance of large-scale polling systems with branching-type and limited service

**Citation for published version (APA):**

Meyfroyt, T. M. M., Boon, M. A. A., Borst, S. C., & Boxma, O. J. (2019). Performance of large-scale polling systems with branching-type and limited service. *Performance Evaluation*, 133, 1-24.  
<https://doi.org/10.1016/j.peva.2019.04.002>

**DOI:**

[10.1016/j.peva.2019.04.002](https://doi.org/10.1016/j.peva.2019.04.002)

**Document status and date:**

Published: 01/09/2019

**Document Version:**

Accepted manuscript including changes made at the peer-review stage

**Please check the document version of this publication:**

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

**General rights**

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

[www.tue.nl/taverne](http://www.tue.nl/taverne)

**Take down policy**

If you believe that this document breaches copyright please contact us at:

[openaccess@tue.nl](mailto:openaccess@tue.nl)

providing details and we will investigate your claim.

# Performance of Large-Scale Polling Systems with Branching-Type and Limited Service

T.M.M. Meyfroyt, M.A.A. Boon, S.C. Borst, O.J. Boxma  
tmeyfroyt@gmail.com, {m.a.a.boon, s.c.borst, o.j.boxma}@tue.nl

Department of Mathematics and Computer Science  
Eindhoven University of Technology, P.O. Box 513  
5600 MB Eindhoven, The Netherlands

May 6, 2019

## Abstract

Motivated by emerging Internet-of-Things (IoT) applications and smart building environments, we analyze the performance of large-scale symmetric polling systems where the number of queues grows large. We consider a scenario in which the total arrival rate is kept fixed and the individual switch-over time and service time distributions remain the same. This asymptotic regime leads to cycles of infinite length and queue lengths with non-trivial distributions. We show that for most traditional service policies the scaled cycle times converge to a deterministic value in the limit, which in turn implies that the queue lengths at the various nodes become asymptotically independent. Using these insights, we find that the behavior of individual queues simplifies to that of a discrete-time bulk service queue in the limit, so that the marginal queue length and waiting-time distributions become considerably easier to analyze. Additionally, we propose a new flexible  $k$ -limited service discipline aimed at striking a good balance between short mean queue lengths and predictable cycle times for deadline-critical applications.

**Keywords:** Polling models, queue lengths, cycle times, flexible  $k$ -limited service.

## 1 Introduction

In the present paper we investigate the performance of large-scale symmetric polling systems. Polling systems provide canonical models for evaluating the delay performance of systems where multiple queues contend for access to a shared resource and are served in an alternating manner. These models find applications in a wide range of domains, e.g., in computer-communications, production and transportation, and many results are available in the literature [3]. However, asymptotic regimes where the number of queues grows large appear to have received hardly any attention so far. Such scenarios are of strong interest in the context of emerging Internet-of-Things (IoT) applications and smart building environments. In building automation and control systems, efficient and economic networking solutions are required in order to transmit and exchange many kinds of monitoring, control, maintenance and management data through the network [13, 22]. Typically the number of devices and sensors in these scenarios is quite large, and it is critical that the end-to-end delay of the data transmitted satisfies predetermined deadlines in order for these systems to meet their performance and functional requirements.

The particular application that motivated the present study is the so-called *BACnet* (Building Automation and Control networks) protocol, which is specifically designed to meet the communication needs of the latter building automation and control systems [8]. BACnet relies on token-passing algorithms for its medium access control. Token-passing algorithms ensure orderly access to a communication channel by passing a token in a cyclic fashion along all nodes in the network. Only when

a node holds the token, it is allowed to send messages to other nodes, before it has to pass on the token to the next node. The token-passing algorithm determines when and to which node the token has to be passed on. Since there is only a single token in the network, only one device can be transmitting at any given time and no data collisions occur. Access to the network is thus guaranteed, and deadline-critical applications can be supported. Another well-known communication protocol that uses token-passing for its medium access control is the Token Ring mechanism [19]. The delay performance of such token-passing algorithms is typically analyzed using polling models.

A polling model is a system that consists of a number of queues that are served by a single server. Since polling systems arise as a useful model in many applications, e.g., in computer-communication, production, transportation and maintenance systems, much work has already been devoted to their analysis [3]. In most cases, it is assumed that customers arrive at the various queues according to independent Poisson processes. Moreover, we assume that the server visits the queues in cyclic order, requiring some switch-over time to move from one queue to the next. At each queue, the server serves a number of customers, which is dictated by the service discipline at that queue. Common service disciplines are exhaustive service (the server keeps serving customers at a queue until there are no more customers at this queue), gated service (the server will serve only those customers that it meets upon its arrival at the queue; other customers arriving to this queue during this visit period will not be served until the next cycle) and  $k$ -limited service (the server serves a queue until either it has become empty or  $k$  customers have been served, whichever event occurs first).

As discussed in the recent survey [5], there is a sharp dichotomy between ‘easy’ and ‘hard’ polling models. If the service discipline at each queue is of so-called branching-type, then a detailed analysis of the joint queue length distribution at all queues is possible via a relation to multi-class branching processes [17] – and waiting-time distributions at all queues can also be obtained. Roughly speaking, the arrivals at the various queues during a service of a customer  $K_i$  at some queue  $Q_i$  may be viewed as the children of  $K_i$ , just as in a branching process. In the case of, e.g., gated service, all  $k$  customers present in a queue at the beginning of a visit of the server to that queue have the same distribution of their numbers of children. Hence, during that visit, all customers present at the beginning of the visit are replaced in a probabilistically identical manner by a new generation of customers. That allows a representation of the numbers of customers in all queues by a multi-class branching process, which can be analyzed exactly. Exhaustive service has a similar property. However, most disciplines, like  $k$ -limited service, do not have this branching property; e.g., under 1-limited, the first customer in the queue is treated differently during a visit than the other customers. If the service discipline in at least one queue is not of branching-type, then only in a few (two-queue) exceptional cases an exact analysis of queue length and waiting-time distributions is known [5].

In view of the usually sizeable number of nodes in BACnet deployments, we focus in the present paper on the performance of polling systems in an asymptotic regime where the number of queues grows large. Many-queue asymptotics of polling systems have been studied before, but only when the mean switch-over times go to zero as the number of queues grows large. In the limit, such a scenario behaves as a “continuous” spatial polling system where a single server moves at a constant rate along a closed tour, stopping to perform services wherever it encounters requests, which appear at locations independently and uniformly distributed over the tour [10, 14]. In this asymptotic regime, the mean cycle time remains finite and queues have either length zero or one with high probability. In contrast, we consider a regime where the number of queues grows large while the individual switch-over times do not shrink to zero in comparison with the service times, but remain the same. This regime gives rise to cycles of infinite length and queue lengths with non-trivial distributions. Performance metrics of interest are the asymptotic queue length, cycle time and waiting-time distributions as the number of queues grows large.

We will first focus on branching-type service disciplines, such as the exhaustive and gated service disciplines. It is well known that the exhaustive service discipline minimizes the total amount of work in the system among all cyclic service policies [16], but the disadvantage is that it does not put any restrictions on the cycle times. For this reason, the exhaustive service discipline is ill-suited for

deadline-critical applications.

A more suitable service discipline for deadline-critical applications is  $k$ -limited service, which will be our second focus. The guarantee that at most  $k$  customers are served at a server visit bounds the cycle time to some extent. This is also the main reason why the BACnet protocol relies on a  $k$ -limited service discipline, since in building automation and control systems it is essential that certain delay-sensitive messages can be sent within a set time frame. However, a major drawback of a  $k$ -limited service discipline is the fact that if the server reaches a queue with less than  $k$  customers, it will not use this unused capacity to serve extra customers in one of the following queues that might contain more than  $k$  customers. In order to overcome this drawback, we introduce a *flexible*  $k$ -limited service discipline, which makes the server more flexible by allowing it to serve more than  $k$  customers in the present queue if less than  $k$  customers were served in one of the previous queues.

Our key contributions are the following. (1) We give explicit results for the covariance of queue lengths, the covariance of visit times and the variance of the cycle time for symmetric polling systems when the server uses a branching-type service discipline. (2) We derive the corresponding many-queue limits. (3) We provide evidence that in the many-queue regime each individual queue behaves asymptotically as a discrete-time bulk service queue, which significantly simplifies the analysis of the marginal queue length and waiting-time distributions for many service disciplines, including the  $k$ -limited service discipline. (4) We introduce the flexible  $k$ -limited service discipline and show how to approximate its performance.

The remainder of this paper is organized as follows. We provide a model description in Section 2. In Section 3 we first derive some exact results for visit times, cycle times and queue lengths in the case of branching-type service disciplines, and we subsequently use those results to obtain many-queue asymptotic results for the cycle time variance and for mean waiting times. We then briefly consider general non-idling policies – service disciplines in which the server does not idle at a queue when that queue has customers – in Section 4, with a focus on the 1-limited service discipline. In Section 5 we propose the flexible  $k$ -limited service policy and provide a way to approximate its asymptotic performance. Finally, in Section 6, we use simulations to compare the performance of the various service disciplines and we investigate how well our asymptotic results can approximate networks of finite size. In Section 7 we summarize our results and discuss the key novel insights and engineering implications.

## 2 Model description and preliminaries

We consider  $n \geq 1$  queues  $Q_1, \dots, Q_n$  being served by a single server who serves at unit speed. The server visits the queues in a cyclic non-idling manner. We assume that customers arrive at the queues according to independent Poisson processes, of rate  $\lambda_i$  at  $Q_i$ ,  $i = 1, \dots, n$ . We assume that customers at  $Q_i$  have i.i.d. service times with first moment  $\beta_i$ , second moment  $\beta_i^{(2)}$  and Laplace-Stieltjes transform (LST)  $\mathcal{B}_i(\cdot)$ ,  $i = 1, \dots, n$ . We define  $\Lambda = \sum_{i=1}^n \lambda_i$ ,  $\rho_i = \lambda_i \beta_i$  and  $\rho = \sum_{i=1}^n \rho_i$ . The switch-over times of the server for moving between  $Q_i$  and the next queue are i.i.d. random variables with first moment  $s_i$ , second moment  $s_i^{(2)}$  and LST  $\mathcal{S}_i(\cdot)$ . All interarrival, service and switch-over times are assumed to be independent.

In addition to the Poisson arrivals, we assume that the service completion of a customer at  $Q_i$  leads to  $M_{i,j}$  additional customers joining  $Q_j$ , where we assume the  $M_{i,j}$  to be random variables which are independent for all  $i$  and  $j$ . This allows that after a customer has received service, with some probability, it stays in the system and is routed to join another queue, which could arise for example in communication networks where some packets are routed in a multi-hop manner. Additionally, it allows that multiple customers simultaneously join the system, after a customer has received service. Again, one can think of communication networks where some packets may require a response from several nodes. By setting  $M_{i,j} = 0$  for all  $i, j$ , our system becomes a standard polling model.

We define the following random variables, for  $i = 1, \dots, n$ :

$V_i^{(k)}$ : The length of the  $k$ 'th server visit to  $Q_i$ .

$S_i^{(k)}$ : The length of the  $k$ 'th switch-over time between  $Q_i$  and  $Q_{i+1}$ ; here and in the sequel,  $S_n^{(k)}$  is the length of a switch-over time from  $Q_n$  to  $Q_1$ .

$C_i^{(k)}$ : The length of the  $k$ 'th cycle *starting* at a visit beginning at  $Q_i$ , i.e.

$$C_i^{(k)} := \sum_{j=i}^n (V_j^{(k)} + S_j^{(k)}) + \sum_{j=1}^{i-1} (V_j^{(k+1)} + S_j^{(k+1)}).$$

$C_i^{*(k)}$ : The length of the  $k$ 'th cycle *starting* at a visit *completion* at  $Q_i$ , i.e.

$$C_i^{*(k)} := S_i^{(k)} + \sum_{j=i+1}^n (V_j^{(k)} + S_j^{(k)}) + \sum_{j=1}^{i-1} (V_j^{(k+1)} + S_j^{(k+1)}) + V_i^{(k+1)}.$$

$I_i^{(k)}$ : The length of the  $k$ 'th intervisit time of  $Q_i$ , i.e.  $I_i^{(k)} := C_i^{(k)} - V_i^{(k)}$ .

Additionally, let  $V_i$ ,  $S_i$ ,  $C_i$ ,  $C_i^*$  and  $I_i$  denote their steady-state limits for  $k \rightarrow \infty$ , assuming these exist. For ease of notation we will write  $V := V_1$ ,  $S := S_1$ ,  $C := C_1$ ,  $C^* := C_1^*$  and  $I := I_1$ .

Lastly, let  $X_j^i$  denote the steady-state queue length of  $Q_j$  at a visit beginning at  $Q_i$  and write  $\mathcal{F}_i(\mathbf{z})$ ,  $\mathbf{z} = (z_1, \dots, z_n)$ , for the probability generating function (PGF) of the steady-state joint queue lengths at a visit beginning at  $Q_i$ .

### 3 Branching-type service disciplines

In this section we will consider polling systems that employ a branching-type service discipline at all queues. After formally introducing the class of branching-type service disciplines, we shall derive the joint transform of the lengths of  $n$  consecutive visit times and switch-over times. That will subsequently give the LST of their sum: the cycle time starting with a visit at, say,  $Q_1$ . The main result of this section is an explicit expression for the cycle time variance, given in Proposition 2, in the case of a symmetric system (all parameters the same for all queues). Subsequently, we use that expression to study the large- $n$  asymptotics of the cycle time variance (cf. Proposition 3).

In Section 1 we argued that polling systems with a branching-type service discipline at all queues can be analyzed in much detail by establishing a direct link to multi-class branching processes. A service discipline is said to be of branching-type when the system satisfies the following property [17]:

**Property 1.** *If there are  $k_i$  customers present at  $Q_i$  at the start of a server visit, then during the course of the visit, each of these  $k_i$  customers will effectively be replaced in an i.i.d. manner by a random population having PGF  $h_i(z_1, \dots, z_n)$ .*

Note that the gated service discipline satisfies Property 1. If we ignore the possibility of having positive numbers  $M_{i,j}$  of additional customers after a service completion, then

$$h_i(z_1, \dots, z_n) = \mathcal{B}_i\left(\sum_{j=1}^n \lambda_j(1 - z_j)\right), \quad (1)$$

the PGF of the numbers of arrivals at all queues during one service time  $B_i$ . The exhaustive service discipline also satisfies Property 1, with

$$h_i(z_1, \dots, z_n) = \theta_i\left(\sum_{j \neq i} \lambda_j(1 - z_j)\right), \quad (2)$$

where  $\theta_i(\cdot)$  is the LST of the time that the server spends at  $Q_i$  due to the presence of one customer there, i.e. the time it spends serving that customer and its descendants during that visit, see [17] and [21] where this is referred to as a sub-busy period. In contrast, the  $k$ -limited service discipline does *not* satisfy Property 1.

Consider now a subclass of polling systems which in addition to Property 1 satisfy the following

**Property 2.**  $M_{i,j} = 0$  for all  $i \neq j$ .

Property 2 essentially puts certain restrictions on the branching functions  $h_i(z_1, \dots, z_n)$ , as it implies the following

$$h_i(z_1, \dots, z_{i-1}, 1, z_{i+1}, \dots, z_n) = \theta_i \left( \sum_{j \neq i} \lambda_j (1 - z_j) \right). \quad (3)$$

For example, this prohibits that a customer, after it has received service, is routed to receive service at a different queue. However, it does permit that a customer immediately returns to the queue where it has just received service. Note that systems using a classical service discipline such as the exhaustive, gated and binomial service disciplines without routing all satisfy Property 2. Binomial-gated is a generalization of the gated discipline: only those customers present in the queue upon the server's arrival (to this queue) are candidates for service during its present visit. Each of these customers will be served with probability  $p$  during this visit. Customers that are skipped stay in the queue and become eligible for service again in the next cycle. The binomial-exhaustive discipline is a similar generalization of the exhaustive discipline; we refer to Subsection 3.4 for a formal definition of binomial-gated and binomial-exhaustive service.

For systems satisfying Properties 1 and 2, we now relate the joint LST of  $V_i + S_i$ ,  $i = 1, \dots, n$ , to the joint PGF  $\mathcal{F}_1(\cdot)$ , extending Corollary 3.1 of [7], which deals with the LST of the cycle time  $C = \sum_{i=1}^n (V_i + S_i)$ . In the following subsections, we will use this result to determine the cycle time variance. In Appendix A, we generalize this result to systems where the restrictive Property 2 is relaxed. However, it is not at all straightforward to extend *all* other results of this section to this more general case.

**Proposition 1.** Write  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_n)$ . In case Properties 1 and 2 hold, the joint LST of  $V_i + S_i$ ,  $i = 1, \dots, n$ , is given by

$$\mathbb{E} \left[ e^{-\sum_{i=1}^n \omega_i (V_i + S_i)} \right] = \mathcal{F}_1(\theta_1(\gamma_1(\boldsymbol{\omega})), \dots, \theta_n(\gamma_n(\boldsymbol{\omega}))) \prod_{i=1}^n \mathcal{S}_i(\gamma_i(\boldsymbol{\omega})), \quad (4)$$

where  $\gamma_n(\boldsymbol{\omega}) = \omega_n$  and, for  $1 \leq i \leq n-1$ ,

$$\gamma_i(\boldsymbol{\omega}) := \omega_i + \sum_{j=i+1}^n \lambda_j (1 - \theta_j(\gamma_j(\boldsymbol{\omega}))). \quad (5)$$

*Proof.* Let  $\mathbf{X} = (X_1^1, \dots, X_n^1)$  and  $\mathbf{m} = (m_1, \dots, m_n)$ . It is sufficient to show that

$$\mathbb{E} \left[ e^{-\sum_{i=1}^n \omega_i (V_i + S_i)} \middle| \mathbf{X} = \mathbf{m} \right] = \prod_{i=1}^n \mathcal{S}_i(\gamma_i(\boldsymbol{\omega})) \theta_i^{m_i}(\gamma_i(\boldsymbol{\omega})).$$

First, notice that, because of Properties 1 and 2, the PGF of the number of customers joining  $Q_j$  during a visit plus switch-over time of  $Q_i$  of length  $x$ , where  $i \neq j$ , is given by

$$\mathbb{E} \left[ z^{X_j^{i+1} - X_j^i} \middle| V_i + S_i = x \right] = e^{-\lambda_j x (1 - z)},$$

which implies that the LST of the time the server spends at  $Q_j$  during its next visit because of these customers is given by  $e^{-\lambda_j x (1 - \theta_j(\omega_j))}$ .

This then allows us to write

$$\begin{aligned}
& \mathbb{E} \left[ e^{-\sum_{i=1}^n \omega_i (V_i + S_i)} \mid \mathbf{X} = \mathbf{m} \right] \\
&= \mathcal{S}_n(\omega_n) \theta_n^{m_n}(\omega_n) \mathbb{E} \left[ e^{-\sum_{i=1}^{n-1} (\omega_i + \lambda_n (1 - \theta_n(\omega_n))) (V_i + S_i)} \mid \mathbf{X} = \mathbf{m} \right] \\
&= \mathcal{S}_n(\omega_n) \theta_n^{m_n}(\omega_n) \mathcal{S}_{n-1}(\gamma_{n-1}(\boldsymbol{\omega})) \theta_{n-1}^{m_{n-1}}(\gamma_{n-1}(\boldsymbol{\omega})) \\
&\quad \cdot \mathbb{E} \left[ e^{-\sum_{i=1}^{n-2} (\omega_i + \lambda_{n-1} (1 - \theta_{n-1}(\gamma_{n-1}(\boldsymbol{\omega}))) + \lambda_n (1 - \theta_n(\omega_n))) (V_i + S_i)} \mid \mathbf{X} = \mathbf{m} \right].
\end{aligned}$$

Performing this procedure  $n$  times and deconditioning gives the desired result.  $\square$

Note that by substituting  $\omega_i = u$  for  $i = 1, \dots, n$  in Equation (4), we readily find the LST of the cycle time  $C$ , as shown by the following corollary.

**Corollary 1.** *The LST of the cycle time  $C = C_1$  is given by*

$$\mathbb{E}[e^{-uC}] = \mathcal{F}_1(\theta_1(\gamma_1(u, \dots, u)), \dots, \theta_n(\gamma_n(u, \dots, u))) \prod_{i=1}^n \mathcal{S}_i(\gamma_i(u, \dots, u)).$$

In the remainder of this paper, we will focus on symmetric systems, i.e.  $\lambda_i = \lambda = \Lambda/n$ ,  $s_i = s$ ,  $s_i^{(2)} = s^{(2)}$ ,  $\beta_i = \beta$  and  $\beta_i^{(2)} = \beta^{(2)}$ ,  $h_i(z_1, \dots, z_n) = h_1(z_i, \dots, z_n, z_1, \dots, z_{i-1})$  and  $\theta_i(u) = \theta(u)$  for all  $i = 1, \dots, n$ . We assume that these second moments,  $s_i^{(2)}$  and  $\beta_i^{(2)}$ , are finite whenever they appear in our formulas. Additionally, we assume that the  $M_{i,j}$  are independent and identically distributed for all  $i$  and  $j \neq i$ , i.e. all queues that are not being served are treated equally. In Subsection 3.2 we will use Proposition 1 to derive  $r_i := \mathbb{E}[V_i + S_i]$  and  $r_{i,j} := \text{Cov}[V_i + S_i, V_j + S_j]$ , and obtain an explicit expression for  $\text{Var}[C]$ . However, in preparation for this, we first determine in Subsection 3.1 explicit expressions for  $\mathbb{E}[X_i^1]$  and  $\text{Cov}[X_i^1, X_j^1]$ .

### 3.1 Queue length covariance

In this subsection we will derive the mean queue lengths  $l_i := \mathbb{E}[X_i^1]$  and the crossmoments  $l_{i,j} := \mathbb{E}[X_i^1 X_j^1]$  for polling systems satisfying Property 1, which will give us  $\text{Var}[X_i^1]$  and  $\text{Cov}[X_i^1, X_j^1]$ .

(i) **Determination of  $l_i$ ,  $i = 1, 2, \dots, n$ .**

Recall that  $\mathcal{F}_i(\mathbf{z})$  is the PGF of the joint queue lengths at a visit beginning at  $Q_i$ . In [17] it is shown that the following relation holds for all polling systems satisfying Property 1:

$$\mathcal{F}_{i+1}(\mathbf{z}) = \mathcal{F}_i(z_1, \dots, z_{i-1}, h_i(\mathbf{z}), z_{i+1}, \dots, z_n) \mathcal{S}_i \left( \sum_{j=1}^n \lambda_j (1 - z_j) \right).$$

For a symmetric system, this gives

$$\mathcal{F}_1(\mathbf{z}) = \mathcal{F}_1(h_1(z_n, z_1, \dots, z_{n-1}), z_1, \dots, z_{n-1}) \mathcal{S} \left( \sum_{j=1}^n \lambda (1 - z_j) \right). \quad (6)$$

We define

$$\left. \frac{\partial}{\partial z_i} h_1(\mathbf{z}) \right|_{\mathbf{z}=(1, \dots, 1)} := \begin{cases} \phi, & i = 1, \\ \psi, & i \neq 1. \end{cases}$$

In the polling literature  $1 - \phi$  is called the *exhaustiveness* of a (branching-type) service discipline, introduced by Van der Mei and Levy [21].

Differentiating (6) with respect to  $z_i$ , we find

$$l_i = \begin{cases} l_{i+1} + \psi l_1 + \lambda s, & i \neq n, \\ \phi l_1 + \lambda s, & i = n. \end{cases}$$

Solving gives

$$l_i = \frac{ns\lambda}{1 - \phi - (n-1)\psi} \left( 1 - \frac{i-1}{n}(1 - \phi + \psi) \right). \quad (7)$$

Note that the stability condition for the system is  $\phi + (n-1)\psi < 1$ , see [17]. For the most common service disciplines satisfying Properties 1 and 2 including (binomial-) gated and (binomial-) exhaustive, we have  $\phi = y + \psi$  and  $\psi = \mu\Lambda/n$  for  $y := \mathbb{E}[M_{1,1}]$  (which is zero for standard gated and exhaustive service) and  $\mu := -\frac{d}{du}\theta(u)|_{u=0}$ . For the systems considered in this paper, neither  $y$  nor  $\mu$  depends on  $n$ , meaning that the total workload arriving to the system per time unit is equal to  $\phi + (n-1)\psi = y + \mu\Lambda$ , which does not depend on  $n$  either. It also follows that  $l_1$  does not depend on  $n$ .

**(ii) Determination of  $l_{i,j}$ ,  $i, j = 1, 2, \dots, n$ .**

Now, define

$$\frac{\partial^2}{\partial z_i \partial z_j} h_1(z) \Big|_{z=(1,\dots,1)} := \begin{cases} \phi^{(2)}, & i = 1, j = 1, \\ \psi_1^{(2)}, & i = j, i \neq 1, \\ \psi_2^{(2)}, & i \neq 1, j \neq 1, i \neq j, \\ \chi, & i = 1, j \neq 1. \end{cases}$$

Differentiating Equation (6) with respect to  $z_i$  and  $z_j$  gives for  $i \neq n, j \neq n, i \neq j$ :

$$l_{i,j} = l_{i+1,j+1} + \psi(l_{1,i+1} + l_{1,j+1}) + \psi^2 l_{1,1} + \lambda s(l_{i+1} + l_{j+1}) + (\psi_2^{(2)} - \psi^2 + 2\lambda s\psi)l_1 + \lambda^2 s^{(2)}; \quad (8)$$

for  $i < n$  and  $j = n$ :

$$l_{i,n} = \phi l_{1,i+1} + \phi \psi l_{1,1} + \lambda s l_{i+1} + (\chi - \phi \psi + \lambda s(\phi + \psi))l_1 + \lambda^2 s^{(2)}; \quad (9)$$

for  $i = j$  and  $i < n$ :

$$l_{i,i} = l_i - l_{i+1} + l_{i+1,i+1} + 2\psi l_{1,i+1} + \psi^2 l_{1,1} + 2\lambda s l_{i+1} + (\psi_1^{(2)} - \psi^2 + 2\lambda s\psi)l_1 + \lambda^2 s^{(2)}; \quad (10)$$

and finally for  $i = j = n$ :

$$l_{n,n} = l_n + \phi^2 l_{1,1} + (\phi^{(2)} - \phi^2 + 2\lambda s\phi)l_1 + \lambda^2 s^{(2)}. \quad (11)$$

Consider Equation (8). By subtracting  $l_{i,j+1}$ , we find for  $j+1 \neq n$  and  $i \neq j+1$ :

$$l_{i,j} - l_{i,j+1} = l_{i+1,j+1} - l_{i+1,j+2} + \psi(l_{1,j+1} - l_{1,j+2}) + \lambda s(l_{j+1} - l_{j+2}).$$

Notice that the last term is independent of  $j$ , suggesting that  $l_{i,j}$  depends linearly on  $j$  for fixed  $i < j$ . Similarly, subtracting  $l_{i+1,j}$  from both sides of Equation (8) gives

$$l_{i,j} - l_{i+1,j} = l_{i+1,j+1} - l_{i+2,j+1} + \psi(l_{1,i+1} - l_{1,i+2}) + \lambda s(l_{i+1} - l_{i+2}),$$

suggesting that  $l_{i,j}$  also depends linearly on  $i$  for fixed  $j > i$ . Together with Equations (9) and (11), assuming  $l_{i,j} = a_i + b_j(i-1) + c_i(j-1)$  then allows us to express  $b_i$  and  $c_i$  in terms of  $b_1, c_1$  and  $l_{1,1}$ . Similarly, rewriting Equation (10), we have

$$l_{i,i} - l_{i+1,i+1} = l_i - l_{i+1} + 2\psi l_{1,i+1} + \psi^2 l_{1,1} + 2\lambda s l_{i+1},$$

which suggests that  $l_{i,i}$  is a quadratic function in  $i$ . Writing  $l_{i,i} = l_{1,1} + b(i-1) + c(i-1)^2$  and substituting then allows us to solve the resulting equations and determine the remaining unknowns, confirming the linear and quadratic relations alluded to above.

**(iii) Determination of  $\text{Var}[X_i^1]$  and  $\text{Cov}[X_i^1, X_j^1]$ .**

We find, after a substantial amount of algebraic manipulation and simplification, that the complete solution is of the form

$$\text{Var}[X_i^1] = l_{i,i} - l_i^2 = \alpha_1 + \alpha_2 - \alpha_3((i-1)(\phi - \psi) - (n-i+1)(1-(i-1)\psi)) - \alpha_4(i-1), \quad (12)$$

and, for  $i < j$ ,

$$\text{Cov}[X_i^1, X_j^1] = l_{i,j} - l_i l_j = \alpha_1 - \alpha_3((i-1)(\phi - \psi) - (n-j+1)(1-(i-1)\psi)), \quad (13)$$

where

$$\begin{aligned} \alpha_1 = & \frac{n\lambda^2(\phi - \psi)(s^{(2)} - s^2)}{(1 + \phi)(1 - \phi - (n-1)\psi)} \\ & + \frac{ns\lambda(\phi - \psi)\psi(\phi^{(2)} + (n-1)\psi_1^{(2)})}{(1 + \phi)(1 - \phi + \psi)(1 - \phi - (n-1)\psi)^2} \\ & + \frac{ns\lambda\chi(1 - \phi^2 - (n-1)\psi^2 - n\psi(1 - \phi))}{(1 + \phi)(1 - \phi + \psi)(1 - \phi - (n-1)\psi)^2} \\ & - \frac{ns\lambda\psi_2^{(2)}(1 - \phi^2 - (n-1)\psi^2 - n\phi(1 - \phi))}{(1 + \phi)(1 - \phi + \psi)(1 - \phi - (n-1)\psi)^2}, \end{aligned} \quad (14)$$

$$\begin{aligned} \alpha_2 = & \frac{ns\lambda(1 + \psi)(1 + \phi^{(2)} - \chi + (n-1)(\psi_1^{(2)} - \psi_2^{(2)}))}{(1 + \phi)(1 - \phi + \psi)(1 - \phi - (n-1)\psi)} \\ & + \frac{ns\lambda\phi(\psi - \phi + \psi_2^{(2)} - \chi)}{(1 + \phi)(1 - \phi + \psi)(1 - \phi - (n-1)\psi)}, \end{aligned} \quad (15)$$

$$\begin{aligned} \alpha_3 = & \frac{ns\lambda(\psi^2\phi^{(2)} + 2(1 - \phi)\psi\chi + (1 - \phi)^2\psi_2^{(2)} + (n-1)\psi^2(\psi_1^{(2)} - \psi_2^{(2)}))}{(1 + \phi)(1 - \phi + \psi)(1 - \phi - (n-1)\psi)^2} \\ & + \frac{\lambda^2(1 - \phi + \psi)(s^{(2)} - s^2)}{(1 + \phi)(1 - \phi - (n-1)\psi)}, \end{aligned} \quad (16)$$

and

$$\alpha_4 = \frac{s\lambda(1 - \phi + \psi + n(\psi_1^{(2)} - \psi_2^{(2)}))}{1 - \phi - (n-1)\psi}.$$

### 3.2 Station-time covariance

We will now use Proposition 1 to derive the mean  $r_i := \mathbb{E}[V_i + S_i]$  and covariance  $r_{i,j} := \text{Cov}[V_i + S_i, V_j + S_j]$  of the station time (sum of the visit time at a queue and switch-over time to the next queue), expressing them in terms of  $\text{Var}[X_i^1]$  and  $\text{Cov}[X_i^1, X_j^1]$ , cf. Equations (12) and (13). This will then allow us to give an explicit expression for  $\text{Var}[C]$  (cf. Proposition 2). Since we will rely on Proposition 1, we restrict ourselves to systems satisfying both Properties 1 and 2 for the remainder of this section.

Define  $\mu^{(2)} := \frac{d^2}{du^2} \theta(u) \Big|_{u=0}$ . First, note that by differentiating (5), we have

$$\Gamma_i^j := \frac{\partial}{\partial \omega_i} \gamma_j(\boldsymbol{\omega}) \Big|_{\boldsymbol{\omega}=(0,\dots,0)} = \mathbb{1}[i=j] + \mu\lambda \sum_{k=j+1}^i \Gamma_i^k.$$

Solving gives

$$\Gamma_i^j = \begin{cases} 1, & i = j, \\ \mu\lambda(1 + \mu\lambda)^{i-j-1}, & j < i, \\ 0, & \text{otherwise.} \end{cases} \quad (17)$$

Differentiating (4) with respect to  $\omega_i$  yields, as expected,

$$r_i = \sum_{j=1}^i (s + \mu l_j) \Gamma_i^j = s + \mu l_1,$$

where the second equality follows from the fact that for service disciplines satisfying Property 2 we have  $\psi = \mu\lambda$ , see Equation (3). Additionally, it follows that

$$\mathbb{E}[C] = nr_i = \frac{ns(1 - \phi + \psi)}{1 - \phi - (n-1)\psi}. \quad (18)$$

Differentiating (4) with respect to  $\omega_1$  and  $\omega_i$  and subtracting  $r_1 r_i = \sum_{j=1}^i \Gamma_i^j (s + \mu l_1)(s + \mu l_j)$  from both sides gives

$$r_{1,i} = \Gamma_i^1 (s^{(2)} - s^2 + l_1(\mu^{(2)} - \mu^2)) + \mu^2 \sum_{j=1}^i \Gamma_i^j (l_{1,j} - l_1 l_j).$$

Using Equation (17), some manipulation yields

$$r_{1,i} = \begin{cases} s^{(2)} - s^2 + l_1(\mu^{(2)} - \mu^2) + \mu^2 \text{Var}[X_1^1], & i = 1, \\ \mu\lambda r_{1,1} + \mu^2 \text{Cov}[X_1^1, X_2^1], & i = 2, \\ (1 + \mu\lambda) r_{1,i-1} - \mu^2 \text{Cov}[X_1^1, X_{i-1}^1 - X_i^1], & i > 2. \end{cases} \quad (19)$$

Equation (13) then tells us that  $\text{Cov}[X_1^1, X_{i-1}^1 - X_i^1] = \alpha_3$  for all  $i > 2$ , hence

$$r_{1,i} = \begin{cases} s^{(2)} - s^2 + \frac{ns\lambda(\mu^{(2)} - \mu^2)}{1 - \phi - (n-1)\psi} + \mu^2(\alpha_1 + \alpha_2 + n\alpha_3), & i = 1, \\ \frac{\mu}{\lambda} \alpha_3 + \frac{\mu}{\lambda} (1 + \mu\lambda)^{i-2} (\lambda^2 r_{1,1} + \alpha_1 \mu\lambda - (1 - (n-1)\mu\lambda)\alpha_3), & i \geq 2. \end{cases} \quad (20)$$

Note that in general it may not hold that  $r_{1,i} > 0$ . As we will see in the next section however,  $r_{1,i}$  will be positive for traditional branching-type service policies, such as the binomial-gated and binomial-exhaustive service disciplines, as  $n$  grows large.

Finally, we find an explicit formula for the variance of the cycle time.

**Proposition 2.**

$$\begin{aligned} \text{Var}[C] &= \sum_{i=1}^n \sum_{j=1}^n r_{i,j} = nr_{1,1} + 2 \sum_{i=1}^{n-1} (n-i) r_{1,1+i} \\ &= 2(1 + \mu\lambda)^n \left( \frac{r_{1,1}}{\mu\lambda} + \frac{\alpha_1}{\lambda^2} - \frac{\alpha_3}{\mu\lambda^3} (1 - (n-1)\mu\lambda) \right) - \frac{2 + n\mu\lambda}{\mu\lambda} r_{1,1} \\ &\quad - \frac{2(1 + n\mu\lambda)}{\lambda^2} \alpha_1 + \frac{2 + \mu\lambda(2 - n(n-1)\mu\lambda)}{\lambda^3 \mu} \alpha_3, \end{aligned} \quad (21)$$

with  $\alpha_1$  given in (14),  $\alpha_3$  in (16), and the  $r_{1,i}$  in (20).

**Remark 1.** For some service disciplines, for example exhaustive service (cf. [2, 23]), it is more natural to consider  $C^*$ , the time between successive visit completions at  $Q_1$ . Obviously it holds that  $\mathbb{E}[C] = \mathbb{E}[C^*]$ , but higher moments generally differ. For the variance, it is readily seen that

$$\text{Var}[C^*] = \text{Var}[S_1 + V_2 + S_2 + \cdots + S_n + V_{n+1}] = \text{Var}[C] + 2 \sum_{i=1}^n \text{Cov}[S_1, V_{i+1}].$$

Following the exact same steps as in the proof of Proposition 1 it is possible to find the joint LST of  $S_1, V_2, S_2, \dots, S_n, V_{n+1}$  which, after differentiation, yields

$$\text{Cov}[S_1, V_{i+1}] = (s^{(2)} - s^2)\lambda\mu(1 + \lambda\mu)^{i-1},$$

for  $i = 1, \dots, n$ . Finally, we obtain

$$\text{Var}[C^*] = \text{Var}[C] + 2(s^{(2)} - s^2)((1 + \lambda\mu)^n - 1), \quad (22)$$

which is strictly greater than  $\text{Var}[C]$  unless  $\text{Var}[S] = 0$ .

### 3.3 Many-queue asymptotics

We will now apply the results of the previous subsections to analyze the limiting behavior of symmetric polling systems satisfying Properties 1 and 2 as the number of queues grows large. To this end, we will consider a sequence of polling systems where the total arrival rate is kept constant, i.e. we now assume  $\lambda = \Lambda/n$  with  $\Lambda$  fixed, so  $\lambda \downarrow 0$  as  $n \rightarrow \infty$ . The main result of this subsection is an explicit expression for the limit of the scaled cycle time variance  $n\text{Var}[C/n]$  (cf. Proposition 3).

First, define the following limiting values

$$\lim_{n \rightarrow \infty} \phi = \Phi, \quad \lim_{n \rightarrow \infty} n\psi = \Psi, \quad \lim_{n \rightarrow \infty} \mu = m, \quad \lim_{n \rightarrow \infty} \mu^{(2)} = m^{(2)}. \quad (23)$$

Additionally, define

$$\lim_{n \rightarrow \infty} \phi^{(2)} = \Phi^{(2)}, \quad \lim_{n \rightarrow \infty} n\psi_1^{(2)} = \Psi_1^{(2)}, \quad \lim_{n \rightarrow \infty} n^2\psi_2^{(2)} = \Psi_2^{(2)}, \quad \lim_{n \rightarrow \infty} n\chi = \mathcal{X}. \quad (24)$$

Note that the limits in (23) should exist by the stability condition  $\phi + (n-1)\psi < 1$ . However, the stability condition does not guarantee the existence of the limits in (24), but their existence is necessary for the existence of a finite limiting value of  $n\text{Var}[C/n]$ .

Additionally, since we limit ourselves to service disciplines satisfying Property 2, we know from Equation (3) that  $\Psi = m\Lambda$  and  $\Psi_1^{(2)} = 0$ .

**Remark 2.** Note that  $\Phi^{(2)}$ ,  $\Psi_2^{(2)}$  and  $\mathcal{X}$  need not be zero, as illustrated by the following example. Consider a system in which the queues get served according to the usual gated service discipline, but whenever a customer completes its service,  $Y$  new customers join the same queue at which the customer received its service, where  $Y$  is a non-negative integer-valued random variable with PGF  $G_Y(z)$ . The branching function is then given by

$$h(\mathbf{z}) = G_Y(z_1) \mathcal{B} \left( \frac{\Lambda}{n} \sum_{i=1}^n (1 - z_i) \right).$$

Consequently, for such a system, we find  $\Phi^{(2)} = \mathbb{E}[Y(Y-1)]$ ,  $\Psi_2^{(2)} = \Lambda^2 m^{(2)}$ , and  $\mathcal{X} = \Lambda m \mathbb{E}[Y]$  with stability condition  $\mathbb{E}[Y] + \Lambda m < 1$ .

We will now derive the limiting value of  $n\text{Var}[C/n]$  as  $n$  grows large. From Equations (18) and (7) we find

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{E}[C/n] &= \frac{s}{1 - \Phi - \Psi}, \\ \lim_{n \rightarrow \infty} l_{[xn]} &= \frac{\Lambda s}{1 - \Phi - \Psi} (1 - x(1 - \Phi)). \end{aligned}$$

Furthermore, write

$$A_1 := \lim_{n \rightarrow \infty} n\alpha_1 = \frac{\Lambda^2 \Phi (s^{(2)} - s^2)}{(1 + \Phi)(1 - \Phi - \Psi)} + \frac{\Lambda s (1 + \Phi - \Psi) \mathcal{X}}{(1 + \Phi)(1 - \Phi - \Psi)^2} + \frac{\Lambda s \Psi \Phi \Phi^{(2)}}{(1 - \Phi^2)(1 - \Phi - \Psi)^2} + \frac{\Lambda s \Phi \Psi_2^{(2)}}{(1 + \Phi)(1 - \Phi - \Psi)^2}, \quad (25)$$

$$A_2 := \lim_{n \rightarrow \infty} \alpha_2 = \frac{\Lambda s}{1 - \Phi - \Psi} \left( 1 + \frac{\Phi^{(2)}}{1 - \Phi^2} \right), \quad (26)$$

$$A_3 := \lim_{n \rightarrow \infty} n^2 \alpha_3 = \frac{\Lambda^2 (1 - \Phi) (s^{(2)} - s^2)}{(1 + \Phi)(1 - \Phi - \Psi)} + \frac{2\Lambda s \Psi \mathcal{X}}{(1 + \Phi)(1 - \Phi - \Psi)^2} + \frac{\Lambda s \Psi^2 \Phi^{(2)}}{(1 - \Phi^2)(1 - \Phi - \Psi)^2} + \frac{2\Lambda s \Psi_2^{(2)} (1 - \Phi)}{(1 + \Phi)(1 - \Phi - \Psi)^2}, \quad (27)$$

and

$$A_4 := \lim_{n \rightarrow \infty} n\alpha_4 = \frac{\Lambda s (1 - \Phi)}{1 - \Phi - \Psi}. \quad (28)$$

Then taking limits in Equations (12) and (13) we find

$$\lim_{n \rightarrow \infty} \text{Var}[X_{[xn]}^1] = A_2 - xA_4,$$

and, for  $x \neq y$ ,

$$\lim_{n \rightarrow \infty} n \text{Cov}[X_{[xn]}^1, X_{[yn]}^1] = A_1 + (1 - x\Phi - y(1 - \Psi x))A_3.$$

Similarly,

$$R_{1,1} := \lim_{n \rightarrow \infty} r_{1,1} = s^{(2)} - s^2 + \frac{s\Lambda}{1 - \Phi - \Psi} (m^{(2)} - m^2) + m^2 A_2, \quad (29)$$

and

$$\lim_{n \rightarrow \infty} nr_{1,[xn]} = \frac{m}{\Lambda} A_3 + \frac{m}{\Lambda} (\Lambda^2 R_{1,1} + m\Lambda A_1 - (1 - m\Lambda)A_3) e^{m\Lambda x}.$$

Finally, we obtain the following result for the scaled variance of the cycle time.

**Proposition 3.**

$$\lim_{n \rightarrow \infty} n \text{Var}[C/n] = \frac{1}{m\Lambda^3} ((2 - m^2 \Lambda^2) A_3 - \Lambda^2 (2 + m\Lambda) R_{1,1} - 2m\Lambda (1 + m\Lambda) A_1) + \frac{2}{m\Lambda^3} (\Lambda^2 R_{1,1} + m\Lambda A_1 - (1 - m\Lambda) A_3) e^{m\Lambda}, \quad (30)$$

where  $A_1$  is given by (25),  $A_3$  by (27) and  $R_{1,1}$  by (29).

**Remark 3.** For  $C^*$ , the cycle time starting at a visit completion, taking  $n \rightarrow \infty$  in Equation (22) gives

$$\lim_{n \rightarrow \infty} (\text{Var}[C^*] - \text{Var}[C]) = 2(s^{(2)} - s^2)(e^{m\Lambda} - 1),$$

from which it follows that

$$\lim_{n \rightarrow \infty} n \text{Var}[C^*/n] = \lim_{n \rightarrow \infty} n \text{Var}[C/n].$$

**Remark 4.** It is noteworthy that for some service disciplines, the mean waiting time can be expressed in terms of the mean residual cycle time. For example, with  $\rho_1 = \Lambda\beta/n$ , we have

$$\mathbb{E}[W_{\text{gated}}] = (1 + \rho_1) \left( \frac{\text{Var}[C]}{2\mathbb{E}[C]} + \frac{\mathbb{E}[C]}{2} \right),$$

$$\mathbb{E}[W_{\text{exhaustive}}] = (1 - \rho_1) \left( \frac{\text{Var}[C^*]}{2\mathbb{E}[C]} + \frac{\mathbb{E}[C]}{2} \right).$$

As  $n \rightarrow \infty$ , we see that in both cases the scaled mean waiting time  $W/n$  converges to a scaled mean residual cycle time:

$$\begin{aligned}\lim_{n \rightarrow \infty} \mathbb{E}[W_{\text{gated}}/n] &= \lim_{n \rightarrow \infty} (1 + \Lambda\beta/n) \left( \frac{\text{Var}[C/n]}{2\mathbb{E}[C/n]} + \frac{\mathbb{E}[C/n]}{2} \right) = \frac{s}{2(1-\rho)}, \\ \lim_{n \rightarrow \infty} \mathbb{E}[W_{\text{exhaustive}}/n] &= \lim_{n \rightarrow \infty} (1 - \Lambda\beta/n) \left( \frac{\text{Var}[C^*/n]}{2\mathbb{E}[C/n]} + \frac{\mathbb{E}[C/n]}{2} \right) = \frac{s}{2(1-\rho)}.\end{aligned}$$

### 3.4 Binomial service

As an example, we will examine the binomial-gated and binomial-exhaustive service disciplines, which are, respectively, generalizations of the gated and exhaustive disciplines. Under the binomial-gated service discipline, when the server finds  $m$  customers present at the start of a visit period of a queue, it will serve  $N \sim \text{Bin}(m, p)$ ,  $0 < p \leq 1$ , of these customers before switching to the next queue ( $p = 1$  corresponds to gated service). Under the binomial-exhaustive service discipline, the server not only serves these  $N$  customers but also any additional customers that arrive to the queue during its visit. The branching functions for the binomial-gated and binomial-exhaustive disciplines at  $Q_1$  are given by

$$h(z) = (1-p)z_1 + p\mathcal{B} \left( \frac{\Lambda}{n} \sum_{i=1}^n (1-z_i) \right)$$

and

$$h(z) = (1-p)z_1 + p\mathcal{P}_n \left( \frac{\Lambda}{n} \sum_{i=2}^n (1-z_i) \right),$$

respectively. Here  $\mathcal{B}(\cdot)$  is the LST of the service time distribution as before and  $\mathcal{P}_n(\cdot)$  the LST of a busy-period distribution in an M/G/1 queue with arrival rate  $\Lambda/n$  and the same service time distribution.

We now make the observation that for both policies we have the following limits:  $m = p\beta$ ,  $m^{(2)} = p\beta^{(2)}$ ,  $\Phi = 1-p$ ,  $\Psi = p\Lambda\beta$ ,  $\Phi^{(2)} = \Psi_1^{(2)} = \mathcal{X} = 0$  and  $\Psi_2^{(2)} = p\Lambda^2\beta^{(2)}$ . Hence, for the second moment of the cycle time, Equation (30) tells us that in the limit for  $n \rightarrow \infty$  we see no distinction between the binomial-gated or binomial-exhaustive service disciplines.

We find, with  $\rho = \Lambda\beta$ ,

$$R_{1,1} = s^{(2)} - s^2 + \frac{\Lambda s}{1-\rho} \beta^{(2)}$$

and

$$\lim_{n \rightarrow \infty} n\text{Var}[C/n] = \frac{p}{(2-p)(1-\rho)} R_{1,1} + \frac{2(1-p)}{p(2-p)\rho} (e^{p\rho} - 1) R_{1,1}. \quad (31)$$

Also note that, using Equation (19), one can show that, for  $x > 0$ ,

$$\lim_{n \rightarrow \infty} nr_{1, \lfloor xn \rfloor} = \frac{1 + (1-p)(1-\rho)e^{p\rho x}}{(2-p)(1-\rho)} \left( p\rho(s^{(2)} - s^2) + \frac{p\rho s \Lambda \beta^{(2)}}{1-\rho} \right).$$

This shows that in the special case  $p = 1$  this limit does not depend on  $x$  and  $\lim_{n \rightarrow \infty} nr_{1,i} = \lim_{n \rightarrow \infty} nr_{1,j}$  for all  $i \neq 1$  and  $j \neq 1$ . For other values of  $p$  the covariance of the visit times increases as a function of  $x$ .

Besides giving an explicit expression for the limiting value of the scaled cycle time variance, Equation (31) has another interesting consequence. It reveals that, since  $\text{Var}[C/n]$  is of order  $1/n$ , the scaled cycle time  $C/n$  converges in probability to a deterministic value equal to  $s/(1-\rho)$ . Therefore, denoting by  $X_i^j(j)$  the queue length of  $Q_i$  at the start of the  $j$ 'th server visit to  $Q_i$ , this suggests that as  $n \rightarrow \infty$ :

$$X_i^j(j+1) \stackrel{d}{=} Y_i(X_i^j(j), j) + A_i(j),$$

where  $Y_i(m_j, j) \sim \text{Bin}(m_j, 1 - p)$  are independent binomially distributed random variables with parameters  $m_j$  and  $1 - p$ , and  $A_i(j)$  are independent Poisson distributed random variables with parameter  $\Lambda s / (1 - \rho)$ ,  $j = 1, 2, \dots$ . In particular, letting  $j \rightarrow \infty$ , we conclude that the steady-state queue length  $X_i^i$  at the start of a visit to  $Q_i$  satisfies the distributional property

$$X_i^i \stackrel{d}{=} Y_i(X_i^i) + A_i, \quad (32)$$

where  $Y_i(m) \sim \text{Bin}(m, 1 - p)$  is a Binomially distributed random variable with parameters  $m$  and  $1 - p$ , and  $A_i$  is a Poisson distributed random variable with parameter  $\Lambda s / (1 - \rho)$ .

Now observe that if  $X$  is Poisson distributed with parameter  $\alpha$  and if, given  $X = x$ , the random variable  $Y$  is binomially distributed with parameters  $x$  and  $q$ , then  $Y$  is Poisson distributed with parameter  $\alpha q$ . Hence, we conclude from (32) that  $X_i^i$  is Poisson distributed with parameter  $\gamma$  satisfying  $\gamma = (1 - p)\gamma + \Lambda s / (1 - \rho)$ , i.e.,  $\gamma = \frac{\Lambda s}{p(1 - \rho)}$ :

$$X_i^i \sim \text{Poi}\left(\frac{\Lambda s}{p(1 - \rho)}\right). \quad (33)$$

In fact, as will also be seen in the following sections, this convergence of the scaled cycle time suggests that we can approximate the steady-state distribution of the queue length at the beginning of a visit for many service disciplines, not necessarily of branching-type, by examining a relation of the form

$$X_i^i \stackrel{d}{=} f(X_i^i) + A_i, \quad (34)$$

where the two terms on the right-hand side are independent, the (possibly random) function  $f(\cdot)$  is determined by the actual service discipline, and  $A_i$  is a Poisson distributed random variable with parameter  $\frac{\Lambda s}{1 - \rho}$ . For example, in the case of exhaustive service  $f(x) = 0$  and in the case of  $k$ -limited service  $f(x) = (x - k)^+$ , as will be seen in Section 4.1.

This means that, in the limit, for many large-scale polling systems the analysis of the marginal queue length distribution simplifies to the analysis of a simpler discrete-time queueing model with i.i.d. arrivals. Often, such a model is much easier to analyze than the pre-limit polling system.

**Remark 5.** *In the case of branching-type service, where each customer is treated equally, one can write  $f(X_i^i) = \sum_{j=1}^{X_i^i} U_j$ , where the  $U_j$  are i.i.d. random variables. For this case, rewriting (34) in terms of PGFs  $G_{X_i^i}(z)$  and  $G_U(z)$  and iterating, gives*

$$G_{X_i^i}(z) = G_{X_i^i}(G_U(z)) e^{\frac{\Lambda s}{1 - \rho}(z - 1)} = \lim_{j \rightarrow \infty} G_{X_i^i}(G_U^{(j)}(z)) e^{\frac{\Lambda s}{1 - \rho} \sum_{k=0}^{j-1} (G_U^{(k)} - 1)},$$

where  $G_U^{(j)}(z) = G_U(G_U^{(j-1)}(z))$  denotes the  $j$ 'th iterate of  $G_U(z)$ . Assuming that  $\mathbb{E}[U] < 1$ , which is necessary for stability of the system, we have that  $\lim_{j \rightarrow \infty} G_U^{(j)}(z) = 1$ , which gives

$$G_{X_i^i}(z) = e^{\sum_{k=0}^{\infty} \frac{\Lambda s}{1 - \rho} (G_U^{(k)}(z) - 1)}.$$

We close this subsection by considering the mean waiting time for binomial-gated service. As shown in [15],

$$\mathbb{E}[W] = \frac{l_{1,1}}{l_1} \frac{n(2 - p) + p\rho}{2\Lambda},$$

which gives the following limit

$$\lim_{n \rightarrow \infty} \mathbb{E}[W/n] = \left(\frac{1}{2} + \frac{1 - p}{p}\right) \frac{s}{1 - \rho}. \quad (35)$$

Indeed, in the limit for  $n \rightarrow \infty$ , the scaled waiting time is given by a residual cycle time plus an additional geometrically distributed number of cycle times. Note that the same limit holds for

binomial-exhaustive service, since the probability that a customer joins a queue while the server is serving that queue, is negligible for large  $n$ .

Examining Equations (31) and (35), we see that when choosing  $p$ , there is actually a trade-off between a small variance of the cycle time and short mean waiting times. For example, consider the case in which both the service times and switch-over times are exponentially distributed with  $\beta = s = 1$  and  $\Lambda = 4/5$ . Then setting  $p = 2/3$  will give a mean waiting time of one full cycle compared to half a cycle for  $p = 1$ . However, it will also reduce the scaled variance of the cycle time from 48 to 28.4.

It is important to note, however, that this increase of the waiting times is of order  $n$  and the decrease of the standard deviation of the unscaled cycle time is of order  $\sqrt{n}$ . While in very large systems this trade-off might not seem so attractive, it could be of importance in somewhat smaller systems with deadline-critical applications, requiring a predictable return time of the server.

## 4 General non-idling service disciplines

We will now briefly consider general non-idling service disciplines, which are not necessarily of the branching-type. The goal of this section is to show that in general  $C/n$  may be expected to converge to a deterministic quantity for many polling systems, even if they do not satisfy Property 1. We present two main results in this section. First, we formulate a precise condition for the aforementioned convergence (cf. Proposition 4). We then use this result to state a conjecture the limiting behavior of a polling model with  $k$ -limited service.

We continue to concentrate on symmetric systems with  $\lambda = \Lambda/n$  and throughout this section we assume that Property 2 holds, i.e.  $M_{i,j} = 0$  for all  $i$  and  $j$ . Let  $L$ ,  $L^b$  and  $L^e$  denote the steady-state queue length at  $Q_1$  at an arbitrary epoch, a visit beginning at  $Q_1$  and a visit completion at  $Q_1$ , respectively. Furthermore, for a random variable  $U$  with CDF  $F(u)$  and finite mean and second moment, let  $U^R$  denote the overshoot of  $U$  with PDF  $(1 - F(u))/\mathbb{E}[U]$  and mean  $\mathbb{E}[U^R] = \mathbb{E}[U^2]/(2\mathbb{E}[U])$ .

The pseudo-conservation law for cyclic service systems [6] gives:

$$\rho\mathbb{E}[W] = \frac{\rho\Lambda\beta^{(2)}}{2(1-\rho)} + \frac{\rho(s^{(2)} - s^2) + \rho ns^2}{2s} + \frac{\rho^2 ns}{2(1-\rho)} - \frac{\rho^2 s}{2(1-\rho)} + n\beta\mathbb{E}[L^e]. \quad (36)$$

Applying Little's law and simplifying shows

$$\mathbb{E}[L - L^e] = \frac{\Lambda}{n} \left( \frac{\rho}{1-\rho} \frac{\beta^{(2)}}{2\beta} + \frac{s^{(2)}}{2s} - \frac{s}{2(1-\rho)} \right) + \frac{\Lambda s}{2(1-\rho)}. \quad (37)$$

Consider now a tagged customer arriving at  $Q_1$  during an intervisit period. We define  $L^I$  to be the queue length of  $Q_1$  at an arbitrary epoch during an intervisit time and let  $\tilde{L}^b$  and  $\tilde{L}^e$  be the queue length of  $Q_1$  at the start and end of the server visit contained in the cycle in which the customer arrives. Conditioning on the length of the intervisit period in which the tagged customer arrived, we find

$$\mathbb{E}[\tilde{L}^e] = \frac{1}{\mathbb{E}[I]} \int_{x=0}^{\infty} x\mathbb{E}[L^e | I = x]d\mathbb{P}[I \leq x] = \frac{\mathbb{E}[L^e I]}{\mathbb{E}[I]} = \frac{\text{Cov}[L^e, I]}{\mathbb{E}[I]} + \mathbb{E}[L^e]. \quad (38)$$

Also note that a similar relationship holds for  $\mathbb{E}[\tilde{L}^b]$ . The PASTA property implies that the expected queue length of  $Q_1$  just before the arrival of the tagged customer equals  $\mathbb{E}[L^I]$ . Considering the change in the queue length of  $Q_1$  since the last visit completion at  $Q_1$ , we can write

$$\mathbb{E}[L^I] = \mathbb{E}[\tilde{L}^e] + \frac{\Lambda}{n}\mathbb{E}[I^R] = \left( \frac{\text{Cov}[L^e, I]}{\mathbb{E}[I]} + \mathbb{E}[L^e] \right) + \frac{\Lambda}{n}\mathbb{E}[I^R]. \quad (39)$$

Moreover, the Fuhrmann-Cooper decomposition [11] tells us that

$$\mathbb{E}[L^I] = \mathbb{E}[L] - \frac{\Lambda}{n} \frac{\Lambda\beta^{(2)}}{2n(1-\rho/n)}. \quad (40)$$

Equating (39) and (40) and multiplying by  $n$ , we find

$$n\mathbb{E}[L - L^e] - \frac{\Lambda}{n} \frac{\Lambda\beta^{(2)}}{2(1-\rho/n)} - \frac{\text{Cov}[L^e, I]}{\mathbb{E}[I/n]} = \frac{\Lambda}{2} \left( \frac{n\text{Var}[I/n]}{\mathbb{E}[I/n]} + \mathbb{E}[I] \right),$$

which, together with Equation (37), gives

$$\lim_{n \rightarrow \infty} n\text{Var}[C/n] = \frac{\Lambda s}{(1-\rho)^2} \beta^{(2)} + \frac{s^{(2)} - s^2}{1-\rho} - \frac{2}{\Lambda} \lim_{n \rightarrow \infty} \text{Cov}[L^e, I]. \quad (41)$$

Hence, we have proved the following proposition.

**Proposition 4.** *If for a symmetric polling system for which Property 2 holds, one has*

$$\lim_{n \rightarrow \infty} \text{Cov}[L^e, I] > -\infty, \text{ then } \lim_{n \rightarrow \infty} n\text{Var}[C/n] < \infty.$$

**Discussion.** In Equation (41) we express the asymptotic variance of the scaled cycle time in terms of  $\text{Cov}[L^e, I]$ . Therefore, if the limit of  $\text{Cov}[L^e, I]$  exists, this implies that  $C/n$  will converge in probability to a deterministic value, as was the case for binomial service disciplines. As an example, consider again the binomial-gated and binomial-exhaustive service disciplines of Section 3.4. Combining Equations (31) and (41) we find

$$\lim_{n \rightarrow \infty} \text{Cov}[L^e, I] = \Lambda \left( \frac{1-p}{(2-p)(1-\rho)} - \frac{1-p}{p(2-p)\rho} (e^{p\rho} - 1) \right) R_{1,1};$$

for these binomial service disciplines,  $\lim_{n \rightarrow \infty} \text{Cov}[L^e, I]$  indeed exists and is positive. For disciplines that are not of branching-type, we have not been able to show that  $\lim_{n \rightarrow \infty} \text{Cov}[L^e, I]$  exists. This appears to be a very challenging problem. That should not be very surprising, in view of the unavailability of exact results for queue lengths and intervisit times for almost all such polling models, when the number of queues is larger than 2. However, if we study (41) in more detail, taking into account that the variance is always nonnegative, we observe that the only way for (41) to become infinite, is when  $\text{Cov}[L^e, I] \rightarrow -\infty$  for  $n \rightarrow \infty$ . This would be very counterintuitive, since the numerical examples in this paper indicate that  $\text{Cov}[L^e, I] \geq 0$  for all service disciplines discussed in Section 6, where equality is only achieved for exhaustive service. For branching-type service disciplines, where each customer is replaced in an i.i.d. manner by its descendants, a long queue at the end of a visit period (i.e., there are many descendants) oftentimes implies that the queue was very long at the beginning of that particular visit period. There are, however, exceptions to this rule. Obviously,  $\text{Cov}[L^e, I] = 0$  for service disciplines where no descendants join the queue being served (as with exhaustive service). It is even possible to (very artificially) construct service disciplines that result in a negative covariance between  $L^e$  and  $I$ . Consider, for example, a service discipline where no customers join the queue being served, except when zero customers have arrived in all the other queues. This is a branching-type service discipline satisfying Property 2, with  $\text{Cov}[L^e, I] < 0$ . However, for (flexible)  $k$ -limited service,  $\text{Cov}[L^e, I] \geq 0$  because  $L_e > 0$  automatically implies that the maximum number of customers has been served in the preceding visit period, which will immediately result in a longer intervisit period due to the large number of arrivals in the other queues.

## 4.1 Limited service

Consider now the well-known  $k$ -limited service discipline. Under this discipline, during a visit, the server keeps serving customers until either  $k$  customers have been served or the queue becomes empty, whichever occurs first.

A major benefit of the  $k$ -limited service discipline is that it in a way bounds the cycle time, which can be of vital importance for deadline-critical applications. However, the  $k$ -limited service discipline does not satisfy Property 1, making it notoriously difficult to analyze.

Indeed, also the evaluation of the limiting value of  $\text{Cov}[L^e, I]$  in Equation (41) is difficult, even for  $k = 1$ , see Appendix B. However, while not giving explicit results, the calculations in Appendix B do suggest that also for the limited service discipline the limit of  $n\text{Var}[C/n]$  exists. Consequently, it is plausible that in symmetric polling systems with a  $k$ -limited service discipline, the scaled cycle time  $C/n$  will again converge in probability to a deterministic value of  $s/(1-\rho)$ , as was also the case for the binomial service discipline.

Based on this result, we conjecture that as  $n$  grows large, the steady-state queue length at the start of a server visit satisfies the following relation, similar to Equations (32) and (34):

$$X_i^i \stackrel{d}{=} (X_i^i - k)^+ + A_i,$$

where  $A_i$  is a Poisson distributed random variable with parameter  $\frac{\Lambda s}{1-\rho}$ . In Section 6.3 we conduct numerical experiments that confirm this relation in the limit. We find that each individual queue asymptotically behaves as an  $M/D/1$  queue with bulk service and fixed capacity as studied in [1].

For this  $M/D/1$  bulk service queue, let  $\pi_i$  be the steady-state probability that the queue length at the start of a visit equals  $i$  and define  $\Pi(z) = \sum_{i=0}^{\infty} \pi_i z^i$  and  $\nu = \Lambda s/(1-\rho)$ . Then in [1] it is shown that

$$\Pi(z) = \frac{(k-\nu)(z-1) \prod_{i=1}^{k-1} (z-z_i)/(1-z_i)}{z^k e^{\nu(1-z)} - 1}, \quad (42)$$

where the  $z_i$  are the zeros of the denominator within the unit circle. Moreover, one can deduce that

$$E[L^b] = \frac{k-(k-\nu)^2}{2(k-\nu)} + \sum_{i=1}^{k-1} \frac{1}{1-z_i}. \quad (43)$$

We conclude that analyzing the marginal queue length and waiting-time distributions for the  $k$ -limited service discipline becomes considerably easier in large polling systems. In Section 6 we will investigate through simulations how well this approach approximates the queue length distribution for finite  $n$ .

Based on the discussion of this section and Proposition 4, we now formulate the following conjecture:

**Conjecture 1.** *If for a symmetric polling system for which Property 2 holds, one has that*

$$\lim_{n \rightarrow \infty} \text{Cov}[L^e, I] > -\infty,$$

*then the steady-state queue length at the start of a server visit will asymptotically behave as an  $M/D/1$  bulk service queue with varying capacity.*

In the next section, we will further build upon this conjecture to analyze the performance of a flexible  $k$ -limited service discipline and discuss what is exactly meant by bulk service with varying capacity.

## 5 Flexible $k$ -limited service

We have already discussed that the  $k$ -limited service discipline achieves a more predictable cycle time compared to the exhaustive and gated service disciplines, making it more suitable for deadline-critical applications. However, it is also known that the waiting times of a system with the  $k$ -limited service discipline can be large compared to the exhaustive and gated service disciplines. This is mostly caused by the fact that if the server reaches a very long queue, it will still serve at most  $k$  customers, even though it possibly did not have to serve any customers at the previously visited queues.

In order to reduce the effect of this drawback, we introduce a *flexible*  $k$ -limited service discipline. This discipline works the same as the  $k$ -limited service discipline, except that when the server serves less than  $k$  customers at a queue, we allow the server to use this ‘lost’ capacity during the visits to the next queues. This has the following benefit. If by chance there is a single abnormally large queue, while the queues before it are almost empty, the server is allowed to spend more time on the abnormally large queue, reducing waiting times.

More specifically, the flexible  $k$ -limited service discipline works as follows. Denote by  $S(i, j)$  the number of served customers during the  $j$ 'th visit of the server to  $Q_{(i-1 \bmod n)+1}$ . Then upon reaching  $Q_i$  for the  $j$ 'th time, the server will serve at most  $K(i, j)$  customers, where

$$K(i, j) = k + \left( \ell k - \sum_{l=1}^{\ell} S(i-l, j - \mathbb{1}[i-l < 1]) \right)^+,$$

with  $k \in \mathbb{N}$  and  $\ell \in \{1, \dots, n-1\}$ . Note that the server will always be allowed to serve at least  $k$  and at most  $(\ell+1)k$  customers during a visit and that for  $\ell=0$  the flexible  $k$ -limited service policy coincides with the traditional  $k$ -limited policy. Furthermore, it is easily seen that during a cycle at most  $(n+\ell)k$  customers will be served: the server then serves  $(\ell+1)k$  customers at some queue and  $k$  customers at the following  $n-1$  queues.

In the remainder of this section we will present an approximate analysis of the queue length distribution under the flexible  $k$ -limited service policy, focusing our attention on large systems with many queues. The reason we resort to an approximate analysis is obvious: even for ordinary  $k$ -limited service, no exact analysis of queue lengths is known except when  $k=1$  and  $n=2$ . We claim that the insight, obtained in previous sections regarding the behavior of the cycle time variance for large  $n$ , will allow us to come up with an accurate queue length approximation. This will be confirmed in Section 6 via numerical experiments.

### 5.1 Performance for large systems

Again let the arrival rate,  $\lambda_i$ , to  $Q_i$  equal  $\Lambda/n$  for all  $i=1, \dots, n$ . Our analysis in the previous sections suggests that  $C/n \rightarrow s/(1-\rho)$  as  $n \rightarrow \infty$ , and therefore the number of customers that join a queue between two consecutive visit beginnings tends to a Poisson distributed random variable with mean  $\Lambda s/(1-\rho)$  as  $n$  grows large. Therefore, we can argue that, under the flexible  $k$ -limited service discipline, each individual queue asymptotically behaves as an  $M/D/1$  queue with bulk service and varying capacity. That is, again denoting by  $L^b(i, j)$  the queue length of  $Q_i$  at the start of the  $j$ 'th server visit, we have as  $n \rightarrow \infty$ :

$$L^b(i, j+1) = (L^b(i, j) - K(i, j))^+ + A(i, j), \quad (44)$$

where the  $A(i, j)$  are i.i.d.  $\text{Poi}(\Lambda s/(1-\rho))$  random variables.

In order to analyze the limiting behavior of Equation (44) as  $j \rightarrow \infty$ , we treat  $\{K(i, j)\}_{j=1}^{\infty}$  as an i.i.d. sequence of random variables. Note that generally these variables will not be strictly independent, since a node having a low  $K$  value in a given cycle is more likely to have a low  $K$  value in the following cycle as well. However, if the sequence  $\{K(i, j)\}_{j=1}^{\infty}$  is treated as i.i.d. with  $p_m = \mathbb{P}[K = m]$

and  $P_m = \mathbb{P}[K \geq m]$ , we can determine the steady-state distribution of the Markov chain given by (44) as shown in [12]. In fact, the numerical examples in Section 6 seems to confirm that, in the limit, the steady-state queue-length distribution of the actual polling model converges to the steady-state distribution of this Markov chain.

Let  $\pi_i$  again denote the steady-state probability that the queue length at the start of a visit equals  $i$  and define  $\Pi(z) = \sum_{i=0}^{\infty} \pi_i z^i$ . Furthermore, write  $\phi_i(z) = \sum_{j=k(\ell+1)-i}^{k(\ell+1)} p_j z^{-j}$ . Writing again  $\nu = \Lambda s / (1 - \rho)$ , we then have, see [12],

$$\Pi(z) = \frac{\sum_{i=0}^{k(\ell+1)-1} \pi_i (P_i - z^i \phi_{k(\ell+1)-i}(z))}{e^{\nu(1-z)} - \phi_{k(\ell+1)}(z)}.$$

Multiplying both the numerator and the denominator by  $z^{k(\ell+1)}$ , we know by Rouché's theorem that, in addition to the zero in  $z = 1$ , the denominator has exactly  $k(\ell + 1) - 1$  zeros within the unit disk denoted by  $z_1, \dots, z_{k(\ell+1)-1}$ . Moreover, since  $\Pi(z)$  is analytic for  $|z| < 1$ , these zeros should also be the zeros of the numerator. Hence,

$$\Pi(z) = \frac{c(z-1) \prod_{i=1}^{k(\ell+1)-1} (z - z_i)}{z^{k(\ell+1)} (e^{\nu(1-z)} - \phi_{k(\ell+1)}(z))}.$$

Using the fact that  $\Pi(1) = 1$ , we find

$$c = (\mathbb{E}[K] - \nu) \prod_{i=1}^{k(\ell+1)-1} \frac{1}{1 - z_i}.$$

Here the probabilities  $p_i$  determining the function  $\phi_{k(\ell+1)}(z)$  and the zeros  $z_i$  remain to be found.

### 5.1.1 The case $\ell = 1$

Consider the simplest case  $\ell = 1$  in which the server is allowed to serve at most  $2k$  customers at a queue. Supported by our analysis in Section 3, we assume asymptotic independence between the queue lengths of two neighboring queues. Since each time the server visits an empty queue it follows that the next queue is allowed to serve up to  $2k$  customers, we should have  $\pi_0 = p_{2k}$ . More generally, we require  $\pi_i = p_{2k-i}$ , for  $i = 0, \dots, k-1$ , and  $p_k = 1 - \sum_{i=0}^{k-1} \pi_i$ . Hence, this tells us that  $q(z) = z^{2k} \phi_{2k}(z)$  is the moment generating function of  $(2k - K)$ . Substituting, we find

$$\Pi(z) = \frac{(\mathbb{E}[K] - \nu)(z-1) \prod_{i=1}^{2k-1} \frac{z - z_i}{1 - z_i}}{z^{2k} e^{\nu(1-z)} - q(z)}. \quad (45)$$

Since the coefficients of the function  $q(z)$  are given in terms of the  $\pi_i$ , we can not directly find the roots  $z_i$  to determine the probabilities  $\pi_i$ . However, it is possible to determine them using the following iterative approach.

We start with initial estimates for  $\pi_i$ ,  $i = 0, \dots, k-1$ , which then also give estimates for the probabilities  $p_i$  and hence the generating function  $q(z)$ . Using these estimates, one can determine the roots  $z_i$  of the denominator of (45), giving a new estimate for the generating function  $\Pi(z)$  and the probabilities  $\pi_i$ . Iterating this process will give estimates for the true values of  $\pi_i$  and  $p_i$ .

Implementing this approach gives good and fast results for reasonable values of  $k$  ( $k < 20$ ) and  $\nu$  not too close to  $k$ , converging after only a few iterations. In the next section we will investigate how well the actual queue length distribution is approximated using this approach.

Finally, note that the simplest case  $k = 1$  allows a more direct treatment. Using the fact that  $\Pi(0) = \pi_0 = p_2$  allows us to find the relation

$$p_2^2 = -(\mathbb{E}[K] - \nu) \frac{z_1}{1 - z_1}. \quad (46)$$

Using Equation (46), we can write  $z_1$  as a function of  $p_2$  as follows

$$z_1(p_2) := \frac{p_2^2}{p_2^2 - (\mathbb{E}[K] - \nu)} = \frac{p_2^2}{p_2^2 - (1 + p_2 - \nu)}.$$

Therefore,  $p_2$  can be determined as the zero of

$$z_1(p_2)^2 e^{\nu(1-z_1(p_2))} - (1-p_2)z_1(p_2) - p_2,$$

also giving  $p_1$ . Additionally, one can deduce that

$$\mathbb{E}[L^b] = \frac{2 - (2 - \nu)^2}{2(1 + p_2 - \nu)} + \frac{1}{1 - z_1}. \quad (47)$$

### 5.1.2 The case $\ell > 1$

Consider now the case  $\ell > 1$ . This case can be treated similarly to the case  $\ell = 1$ , but it does require solving some combinatorial problems. For example, consider the case  $\ell = 2$  and  $k = 1$ . Let  $X_{-3}, X_{-2}$  and  $X_{-1}$  denote the number of customers the server found at the last three queues that it visited and let  $K_{-3}, K_{-2}$  and  $K_{-1}$  denote the  $K$  values at those queues during those visits. If we again assume that queue lengths of neighboring queues are independent as  $n \rightarrow \infty$ , then the only way the current queue can have  $K = 3$  is if its two predecessors did not have any customers waiting, hence  $p_3 = \mathbb{P}[X_{-2} = X_{-1} = 0] = \pi_0^2$ . Similarly, the probability that a queue has  $K = 2$  is given by

$$\begin{aligned} p_2 &= \mathbb{P}[X_{-2} + X_{-1} = 1] + \mathbb{P}[X_3 > 1, K_{-3} > 1, X_{-2} = 0, X_{-1} > 1] \\ &\quad + \mathbb{P}[X_{-2} > 1, K_{-2} = 1, X_{-1} = 0] \\ &= 2\pi_0\pi_1 + (1 - \pi_0 - \pi_1)(1 - p_1)\pi_0(1 - \pi_0 - \pi_1) + (1 - \pi_0 + \pi_1)p_1\pi_0. \end{aligned}$$

One can imagine that as  $\ell$  and  $k$  become larger, expressing the probabilities  $p_i$  in terms of the  $\pi_j$  becomes more difficult. However, in practice, using moderate values of  $\ell$  should suffice.

Note that in order to avoid dependencies between queue lengths and  $K$ 's of neighboring queues and to simplify the analysis, one could also consider a randomized flexible  $k$ -limited service policy. By randomized we mean that instead of looking at how many customers were served during the last  $\ell$  visits, the service limit  $K(i, j)$  is determined by how many customers were served at  $\ell$  random visits during the last cycle. Specifically, let  $X^{i,j} = \{X_1^{i,j}, \dots, X_\ell^{i,j}\}$  be i.i.d. randomly chosen subsets of  $\{1, \dots, n\} \setminus \{i\}$ , then

$$K(i, j) = k + \left( \ell k - \sum_{l=1}^{\ell} S(X_l^{i,j}, j - \mathbb{1}[X_l^{i,j} > i]) \right)^+.$$

Such a policy admits an easier analysis, since the probabilities  $p_i$  can be expressed more easily in terms of the  $\pi_j$ . Moreover, there should be weaker dependencies between  $K(i, j)$  and  $K(i, j + 1)$ .

## 6 Numerical results

This section contains the results of extensive numerical experiments. In Subsection 6.1 we evaluate how well the limiting many-queue results of the previous sections for the marginal queue length distribution can serve as approximations for polling systems of finite size. We focus on the binomial-gated,  $k$ -limited and flexible  $k$ -limited service disciplines. In Subsection 6.2 we compare these disciplines with respect to the mean queue length and the cycle time variance. Finally, in Subsection 6.3, we numerically study the limiting behavior of  $\text{Cov}[L^\ell, I]$  for these disciplines.

## 6.1 Marginal queue length distribution convergence

We first consider the binomial-gated service discipline by investigating how well Equation (33) approximates the steady-state queue length at the start of a server visit for finite  $n$ . To this end, we simulate a polling system with  $n = 10, 20, 50, 100, 200, 500$  for  $10^8$  cycles. The results for binomial-gated service are obtained by an exact analysis, using the PGF and LST inversion algorithms by Choudhury and Whitt [9]. We assume that service and switch-over times are exponentially distributed with mean  $2/3$  and  $1$ , respectively. Additionally, we set  $\Lambda = 1$ .

$n$	$p = 1/4$	$p = 1/2$	$p = 3/4$	$p = 1$
10	$5.3 \times 10^{-4}$	$1.0 \times 10^{-3}$	$1.7 \times 10^{-3}$	$3. \times 10^{-3}$
20	$1.5 \times 10^{-4}$	$2.9 \times 10^{-4}$	$5.1 \times 10^{-4}$	$9.2 \times 10^{-4}$
50	$2.6 \times 10^{-5}$	$5.1 \times 10^{-5}$	$9.2 \times 10^{-5}$	$1.7 \times 10^{-4}$
100	$6.7 \times 10^{-6}$	$1.3 \times 10^{-5}$	$2.4 \times 10^{-5}$	$4.4 \times 10^{-5}$
200	$1.7 \times 10^{-6}$	$3.4 \times 10^{-6}$	$6.1 \times 10^{-6}$	$1.1 \times 10^{-5}$
500	$2.7 \times 10^{-7}$	$5.4 \times 10^{-7}$	$9.9 \times 10^{-7}$	$1.9 \times 10^{-6}$

Table 1: Binomial-gated service: Squared error between limiting and exact queue length distributions.

In Table 1 we show the squared error between the PDF of the limiting queue length distribution at cycle beginnings, given by Equation (33), and the corresponding exact results. The squared error between two PDFs  $f(\cdot)$  and  $g(\cdot)$  is defined as  $\sum_{i=0}^{\infty} (f(i) - g(i))^2$ . We find that even for small  $n$  the limiting distribution given by Equation (33) already approximates the steady-state queue length distribution well. In fact, as we can expect from the analysis of Section 3, the squared error is of order  $1/n^2$ . Additionally, we see that as  $p$  becomes smaller, the approximation tends to become more accurate.

$n$	$k = 4$	$k = 5$	$k = 6$
10	$3.18 \times 10^{-3}$	$3.88 \times 10^{-3}$	$4.60 \times 10^{-3}$
20	$9.42 \times 10^{-4}$	$1.13 \times 10^{-3}$	$1.33 \times 10^{-3}$
50	$1.70 \times 10^{-4}$	$2.01 \times 10^{-4}$	$2.37 \times 10^{-4}$
100	$4.46 \times 10^{-5}$	$5.25 \times 10^{-5}$	$6.15 \times 10^{-5}$
200	$1.21 \times 10^{-5}$	$1.38 \times 10^{-5}$	$1.57 \times 10^{-5}$
500	$2.02 \times 10^{-6}$	$2.26 \times 10^{-6}$	$2.57 \times 10^{-6}$

Table 2:  $k$ -limited service: Squared error between limiting and simulated queue length distributions.

Consider now the traditional  $k$ -limited service discipline without flexibility. In Table 2 we show the squared error between the PDF of the limiting queue length defined by the PGF of Equation (42) and the corresponding simulation results for several values of  $k$ , where we have used the same parameter settings as before. Also here we find fast convergence to the limiting queue length distribution. For this case, however, the choice of  $k$  does not appear to have a big influence on the convergence rate.

Finally, we consider the flexible  $k$ -limited service discipline with  $\ell = 1$ . We compare simulation results with the PDF of the limiting queue length defined by the PGF of Equation (45), which we determine by the iterative procedure described in Section 5.1. Results can be found in Table 3. As before, we find fast convergence to the limiting queue length distribution. To get a better impression of how well the actual queue length distribution is approximated by the limiting distribution, we have created a plot to compare them. Figure 1 shows the queue length distributions for systems with  $n = 10$  and  $n = 50$  queues, respectively, with flexible  $k$ -limited service ( $k = 4$  and  $\ell = 1$ ). Comparing these probabilities with the limiting distribution, also visualized in Figure 1, we see that

$n$	$k = 4$	$k = 5$	$k = 6$
10	$3.62 \times 10^{-3}$	$2.93 \times 10^{-3}$	$2.91 \times 10^{-3}$
20	$1.13 \times 10^{-3}$	$8.96 \times 10^{-4}$	$8.97 \times 10^{-4}$
50	$2.22 \times 10^{-4}$	$1.66 \times 10^{-4}$	$1.66 \times 10^{-4}$
100	$6.56 \times 10^{-5}$	$4.41 \times 10^{-5}$	$4.34 \times 10^{-5}$
200	$2.00 \times 10^{-5}$	$1.14 \times 10^{-5}$	$1.12 \times 10^{-5}$
500	$4.95 \times 10^{-6}$	$1.74 \times 10^{-6}$	$1.93 \times 10^{-6}$

Table 3: Flexible  $k$ -limited service with  $\ell = 1$ : Squared error between limiting and simulated queue length distributions.

the general shape is quite similar for  $n = 10$ , but the individual probabilities are not very close yet. For  $n = 50$  the differences have almost disappeared and the limiting probabilities would be a very good approximation for the true distribution.

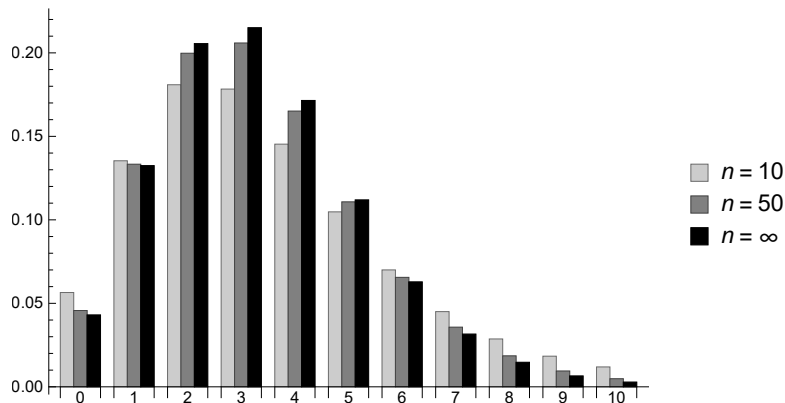


Figure 1: Queue-length distributions for flexible  $k$ -limited service with  $k = 4$  and  $\ell = 1$ .

## 6.2 Comparison of service disciplines

In the previous subsection, we used simulations to examine how well asymptotic results can be used to approximate the performance of polling systems of finite size. In this subsection, we will compare the performance of the binomial-gated,  $k$ -limited and flexible  $k$ -limited service discipline.

We focus our attention on the cycle time variance and the mean queue length. In practice, one seeks a service discipline that minimizes the mean queue length, as this will also have a favorable effect on the mean waiting time of customers. However, as we shall see, lowering the mean queue length usually comes at the price of increased cycle time variance. This trade-off is particularly important for deadline-critical applications, where upon returning to a queue, the server might find a high-priority customer that needs to be served immediately. In such a system a predictable return time of the server is necessary in order to guarantee that these high-priority customers with high probability do not have to wait longer than some predefined threshold.

In Tables 4 and 5 we show the variance of the cycle time and the mean queue length for the binomial-gated service discipline obtained by an exact analysis. The last row of both tables corresponds to the exact asymptotic results from Section 3.4. We find that, as predicted by Equation (31), we can actually decrease the cycle time variance by decreasing  $p$ . As expected (see the remark below Equation (7)), Table 5 shows that the *mean* queue length does not depend on the number of queues.

Consider now the  $k$ -limited and the flexible  $k$ -limited service disciplines. In Tables 6 and 7 we show estimates of the variance of the cycle time and the mean queue length for both disciplines.

$n$	$p = 1/4$	$p = 1/2$	$p = 3/4$	$p = 1$
10	4.91	6.37	8.12	10.3
20	4.95	6.47	8.31	10.6
50	4.98	6.53	8.42	10.9
100	4.98	6.55	8.46	10.9
200	4.99	6.56	8.48	11.0
500	4.99	6.56	8.50	11.0
$\infty$	4.99	6.57	8.50	11.0

Table 4: Binomial-gated service: Scaled variance of the cycle time distribution.

$n$	$p = 1/4$	$p = 1/2$	$p = 3/4$	$p = 1$
10	12	6	4	3
$\infty$	12	6	4	3

Table 5: Binomial-gated service: Mean queue length at the start of a server visit.

$n$	$\ell = 0$			$\ell = 1$		
	$k = 4$	$k = 5$	$k = 6$	$k = 4$	$k = 5$	$k = 6$
10	5.03	7.37	9.01	5.82	8.05	9.27
20	5.19	7.70	9.34	6.26	8.75	9.93
50	5.30	7.93	9.57	6.61	9.28	10.38
100	5.33	8.02	9.65	6.75	9.49	10.54
200	5.35	8.07	9.69	6.82	9.60	10.62
500	5.36	8.10	9.72	6.86	9.68	10.67

Table 6:  $k$ -limited and flexible  $k$ -limited service: Scaled variance of the simulated cycle time distribution.

$n$	$\ell = 0$			$\ell = 1$		
	$k = 4$	$k = 5$	$k = 6$	$k = 4$	$k = 5$	$k = 6$
10	4.21	3.25	2.98	3.72	3.19	3.06
20	4.00	3.22	3.02	3.49	3.11	3.03
50	3.88	3.20	3.04	3.36	3.07	3.01
100	3.84	3.20	3.05	3.31	3.06	3.01
200	3.82	3.20	3.05	3.29	3.05	3.01
500	3.81	3.20	3.06	3.28	3.05	3.01
$\infty$	3.80	3.20	3.06	3.26	3.04	3.01

Table 7:  $k$ -limited and flexible  $k$ -limited service: Simulated mean queue length at the start of a server visit.

Recall that the setting  $\ell = 0$  corresponds to the  $k$ -limited service discipline without flexibility. Comparing with the results in Tables 4 and 5, we find that the  $k$ -limited service discipline is better at reducing the cycle time variance than the binomial-gated service discipline, without increasing the mean queue length too much when compared with the gated service discipline, i.e.  $p = 1$ .

If we now compare the results for the flexible  $k$ -limited service discipline with  $k = 4$  with the results for the  $k$ -limited service discipline without flexibility and the gated service discipline, we find that the flexible  $k$ -limited service discipline is able to achieve the best of both worlds. Compared to the gated service discipline, it achieves almost a 40% decrease in scaled cycle time variance, while

only increasing the mean queue length by roughly 10%. Compared to the  $k$ -limited service discipline without flexibility, it reduces the mean queue length by 15%, while only increasing the scaled cycle time variance by approximately 25%.

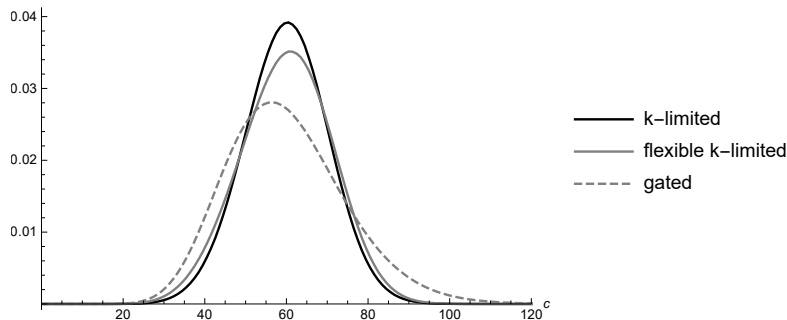


Figure 2: Cycle time PDFs for the  $k$ -limited, flexible  $k$ -limited and gated service disciplines for  $n = 20$  and  $k = 4$ .

In Figure 2 we visualize this trade-off by plotting the cycle time PDFs of the gated service discipline and (simulations for) the  $k$ -limited service discipline with and without flexibility and  $k = 4$ . Indeed, the flexible  $k$ -limited service discipline is able to guarantee a server return time smaller than 85 time units for 99% of the cycles, while the gated service discipline can only guarantee this for 95% of the cycles. Finally, we see that the distributions for the  $k$ -limited service disciplines with and without flexibility are very similar.

### 6.3 Covariance between $L^e$ and $I$

Conjecture 1 in Section 4 states that the covariance between an arbitrary intervisit time and the queue length at the start of this intervisit time remains bounded in the limit as  $n$  tends to infinity. For the (binomial) gated and (binomial) exhaustive service disciplines we have been able to find an exact expression for  $\text{Cov}[L^e, I]$  and prove that its limit exists when  $n \rightarrow \infty$ . For other branching-type service disciplines one can follow a similar approach to find  $\text{Cov}[L^e, I]$  and study its limit when  $n \rightarrow \infty$ . For non-branching service disciplines, we have to resort to simulation to validate the conjecture. Table 8 shows simulation results for  $k$ -limited and flexible  $k$ -limited service. It is clearly visible that  $\lim_{n \rightarrow \infty} \text{Cov}[L^e, I]$  remains bounded as  $n$  grows larger and, as expected, is strictly positive. With 250 simulation runs of  $10^8$  cycles each, the results are accurate up to the two decimals displayed in Table 8 and strongly seem to confirm Conjecture 1.

$n$	$\ell = 0$			$\ell = 1$		
	$k = 4$	$k = 5$	$k = 6$	$k = 4$	$k = 5$	$k = 6$
10	2.75	1.77	1.08	2.42	1.48	0.98
20	2.79	1.64	0.89	2.31	1.16	0.62
50	2.81	1.53	0.75	2.18	0.88	0.35
100	2.81	1.49	0.69	2.12	0.77	0.25
200	2.82	1.47	0.66	2.09	0.70	0.20
500	2.82	1.46	0.64	2.07	0.66	0.17

Table 8:  $k$ -limited and flexible  $k$ -limited service: Simulated covariance between  $L^e$  and  $I$ .

## 7 Conclusion

In the present paper we have provided a detailed analysis of the variance of the cycle time, the covariance of visit times and the covariance of queue lengths in symmetric polling systems where the total number of queues grows large. In contrast to “continuous” polling models, the distribution of the individual switch-over times is assumed to remain the same as the number of queues grows large. While these results are of theoretical merit in their own right, they also yield significant novel insights from a practical perspective that have important engineering implications. Indeed, multi-queue single-server models have found ubiquitous use in a wide range of engineering applications, especially in the context of computer systems and communication networks. The complexity, size and required performance standards in these systems continue to increase due to advances in technology. For example, computing resources are shared among large numbers of users with increasingly tight delay budgets in cloud platforms and virtualized network environments. Wireless communication bandwidth is also shared by increasingly large numbers of users and devices with quite demanding latency and reliability requirements, especially in (tactical) Internet-of-Things applications and smart control environments. As these two examples illustrate, in many of these scenarios the number of contending users is large, with each individual user only requiring a relatively small portion of the overall capacity, but possibly involving quite stringent delay requirements.

The analysis and optimization of such multi-queue models tends to be a major challenge due to the complex characteristics, particularly the dynamic interactions and intricate dependencies that arise among the various queues that contend for service. While polling systems with branching-type service disciplines allow for an exact analysis of the joint queue length distribution as mentioned earlier, these results unfortunately do not immediately offer a recipe for optimization. These issues are further exacerbated for non-branching service disciplines, such as  $k$ -limited or time-limited policies, which are widely used in practical systems, but do not yield an exact analysis, rendering optimization of system parameters even more challenging.

In the present paper we demonstrated that the analysis and optimization of performance in terms of queue lengths and delays dramatically simplifies in symmetric polling systems with a large number of queues. Specifically, the fact that the variance of the scaled cycle time vanishes as the number of queues tends to infinity, means that the scaled cycle time converges to a deterministic quantity. This in turn implies that the various queues behave in a nearly independent manner, while the individual behavior simplifies to that of a stand-alone discrete-time bulk service queue. We have also built on these insights to propose and analyze a new flexible  $k$ -limited service discipline aimed at achieving a good trade-off between short mean queue lengths and predictable cycle times.

Numerical experiments indicate that the discrete-time bulk service queue in fact provides a surprisingly accurate approximation even for a fairly moderate number of queues. We further conjecture that this extends to systems that are not strictly symmetric as long as the load of each individual queue vanishes, in the sense that in the limit each queue behaves as a bulk service queue. The parameters of the latter model only depend on the arrival rate and service time distribution of the queue under consideration, and the aggregate load and total mean switch-over time. Proving a rigorous result along these lines remains as an interesting topic for further research. These novel findings indicate that a complicated high-dimensional optimization problem, such as the selection of service limits, can be decoupled into a collection of one-dimensional optimization problems. This property can then be leveraged to derive simple rules-of-thumb for engineering purposes. In particular, our results provide insight in the delay performance of token-passing algorithms such as the BACnet protocol in deployment scenarios with a large number of nodes and deadline-critical applications which require short delays and guaranteed token return times.

## Acknowledgments

The authors thank Dee Denteneer (Philips Research) for fruitful discussions and useful comments on drafts of the present paper. The research of Onno Boxma and Sem Borst is partly funded by the NWO Gravitation Programme NETWORKS (Grant Number 024.002.003).

## Appendix

### A Extension of Proposition 1

In this appendix, we generalize Proposition 1 to find the cycle time for branching-type polling systems that do *not* satisfy Property 2. Denote by  $C_{i,j}(z)$  the PGF of  $M_{i,j}$  as introduced in Section 2. For simplicity, we assume that the  $M_{i,j}$  are independent, but at the end of this section we briefly discuss the case with dependencies. A visit period at  $Q_i$  starting with  $m$  customers present at  $Q_i$  is considered to consist of  $m$  sub-busy periods, each having LST  $\theta_i(\cdot)$ . Additionally, let  $\tilde{\theta}_i(\mathbf{z}, \boldsymbol{\omega})$  denote the joint transform of the number of customers that are served during a sub-busy period visit to  $Q_i$  and the length of the sub-busy period. Lastly, we define  $L_i$  as the number of service completions at  $Q_i$  during a visit in steady state.

**Proposition 5.** Write  $\mathbf{z} = (z_1, \dots, z_n)$  and  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_n)$ . In case Property 1 holds, the joint transform of  $L_i$  and  $V_j + S_j$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, n$ , is given by

$$\mathbb{E} \left[ \left( \prod_{i=1}^n z_i^{L_i} \right) e^{-\sum_{j=1}^n \omega_j (V_j + S_j)} \right] = \mathcal{F}_1(\tilde{\kappa}_1(\mathbf{z}, \boldsymbol{\omega}), \dots, \tilde{\kappa}_n(\mathbf{z}, \boldsymbol{\omega})) \prod_{i=1}^n \mathcal{S}_i(\tilde{\gamma}_i(\mathbf{z}, \boldsymbol{\omega})), \quad (48)$$

where  $\tilde{\gamma}_n(\mathbf{z}, \boldsymbol{\omega}) = \omega_n$ ,  $\tilde{\kappa}_n(\mathbf{z}, \boldsymbol{\omega}) = \tilde{\theta}_n(z_n, \omega_n)$  and, for  $1 \leq i \leq n-1$ ,

$$\tilde{\gamma}_i(\mathbf{z}, \boldsymbol{\omega}) = \omega_i + \sum_{j=i+1}^n \lambda_j (1 - \tilde{\theta}_j(z_j, \tilde{\gamma}_j(\mathbf{z}, \boldsymbol{\omega}))), \quad (49)$$

and

$$\tilde{\kappa}_i(\mathbf{z}, \boldsymbol{\omega}) = \tilde{\theta}_i \left( z_i \prod_{j=i+1}^n C_{i,j}(\tilde{\kappa}_j(\mathbf{z}, \boldsymbol{\omega})), \tilde{\gamma}_i(\mathbf{z}, \boldsymbol{\omega}) \right). \quad (50)$$

*Proof.* We define

$Z_i :=$  extra work introduced at  $Q_i$  during  $V_1 + S_1$ ,

$H_i :=$  extra customers introduced at  $Q_i$  during  $V_1 + S_1$ .

Consider first the 2-queue case. Then

$$\begin{aligned} & \mathbb{E} \left[ z_1^{L_1} z_2^{L_2} \exp(-\omega_1(V_1 + S_1) - \omega_2(V_2 + S_2)) \mid \mathbf{X} = (m_1, m_2) \right] \\ &= \tilde{\theta}_2(z_2, \omega_2)^{m_2} \mathcal{S}_2(\omega_2) \mathbb{E} \left[ z_1^{L_1} z_2^{H_2} \exp(-\omega_1(V_1 + S_1) - \omega_2 Z_2) \mid \mathbf{X} = (m_1, m_2) \right] \\ &= \tilde{\theta}_1 \left( z_1 C_{1,2}(\tilde{\theta}_2(z_2, \omega_2)), \omega_1 + \lambda_2(1 - \tilde{\theta}_2(z_2, \omega_2)) \right)^{m_1} \tilde{\theta}_2(z_2, \omega_2)^{m_2} \\ & \quad \cdot \mathcal{S}_1(\omega_1 + \lambda_2(1 - \tilde{\theta}_2(z_2, \omega_2))) \mathcal{S}_2(\omega_2), \end{aligned}$$

where for the second equality, we have used that

$$\begin{aligned} & \mathbb{E} \left[ z_2^{H_2} e^{-\omega_2 Z_2} \mid L_1 = l \right] \\ &= \tilde{\theta}_1 \left( C_{1,2}(\tilde{\theta}_2(z_2, \omega_2)), \lambda_2(1 - \tilde{\theta}_2(z_2, \omega_2)) \right)^l \mathcal{S}_1(\lambda_2(1 - \tilde{\theta}_2(z_2, \omega_2))). \end{aligned}$$

It follows that

$$\begin{aligned} & \mathbb{E} \left[ z_1^{L_1} z_2^{L_2} e^{-\omega_1(V_1+S_1)-\omega_2(V_2+S_2)} \right] \\ &= \mathcal{F}_1 \left( \tilde{\theta}_1 \left( z_1 C_{1,2} \left( \tilde{\theta}_2(z_2, \omega_2) \right), \omega_1 + \lambda_2(1 - \tilde{\theta}_2(z_2, \omega_2)) \right), \tilde{\theta}_2(z_2, \omega_2) \right) \\ & \quad \cdot \mathcal{S}_1(\omega_1 + \lambda_2(1 - \tilde{\theta}_2(z_2, \omega_2))) \mathcal{S}_2(\omega_2). \end{aligned}$$

Consider now  $n$  queues. Generalizing the analysis of the 2-queue example, we deduce

$$\begin{aligned} & \mathbb{E} \left[ \left( \prod_{i=1}^n z_i^{L_i} \right) e^{-\sum_{j=1}^n \omega_j(V_j+S_j)} \mid \mathbf{X} = \mathbf{m} \right] \\ &= \mathcal{S}_n(\tilde{\gamma}_n(\mathbf{z}, \boldsymbol{\omega})) (\tilde{\kappa}_n(z_n, \omega_n))^{m_n} \\ & \quad \cdot \mathbb{E} \left[ \left( \prod_{i=1}^{n-1} (z_i C_{i,n}(\tilde{\kappa}_n(\mathbf{z}, \boldsymbol{\omega})))^{L_i} \right) e^{-\sum_{i=1}^{n-1} (\omega_i + \lambda_n(1 - \tilde{\theta}_n(z_n, \omega_n)))(V_i+S_i)} \mid \mathbf{X} = \mathbf{m} \right] \\ &= \mathcal{S}_n(\tilde{\gamma}_n(\mathbf{z}, \boldsymbol{\omega})) (\tilde{\kappa}_n(\mathbf{z}, \boldsymbol{\omega}))^{m_n} \mathcal{S}_{n-1}(\tilde{\gamma}_{n-1}(\mathbf{z}, \boldsymbol{\omega})) (\tilde{\kappa}_{n-1}(\mathbf{z}, \boldsymbol{\omega}))^{m_{n-1}} \\ & \quad \cdot \mathbb{E} \left[ \left( \prod_{i=1}^{n-2} (z_i C_{i,n-1}(\tilde{\kappa}_{n-1}(\mathbf{z}, \boldsymbol{\omega})) C_{i,n}(\tilde{\kappa}_n(\mathbf{z}, \boldsymbol{\omega})))^{L_i} \right) \right. \\ & \quad \left. \cdot e^{-\sum_{i=1}^{n-2} (\omega_i + \lambda_{n-1}(1 - \tilde{\theta}_{n-1}(z_{n-1}, \tilde{\gamma}_{n-1}(\mathbf{z}, \boldsymbol{\omega})) + \lambda_n(1 - \tilde{\theta}_n(z_n, \omega_n)))(V_i+S_i)} \mid \mathbf{X} = \mathbf{m} \right]. \end{aligned}$$

Iterating, we find

$$\mathbb{E} \left[ \left( \prod_{i=1}^n z_i^{L_i} \right) e^{-\sum_{j=1}^n \omega_j(V_j+S_j)} \mid \mathbf{X} = \mathbf{m} \right] = \prod_{i=1}^n \tilde{\kappa}_i(\mathbf{z}, \boldsymbol{\omega})^{m_i} \mathcal{S}_i(\tilde{\gamma}_i(\mathbf{z}, \boldsymbol{\omega})).$$

Consequently,

$$\mathbb{E} \left[ \left( \prod_{i=1}^n z_i^{L_i} \right) e^{-\sum_{j=1}^n \omega_j(V_j+S_j)} \right] = \mathcal{F}_1(\tilde{\kappa}_1(\mathbf{z}, \boldsymbol{\omega}), \dots, \tilde{\kappa}_n(\mathbf{z}, \boldsymbol{\omega})) \prod_{i=1}^n \mathcal{S}_i(\tilde{\gamma}_i(\mathbf{z}, \boldsymbol{\omega})).$$

□

In particular, this implies that the LST of the cycle time is given by

$$\mathbb{E} [e^{-uC}] = \mathcal{F}_1(\kappa_1(u), \dots, \kappa_n(u)) \prod_{i=1}^n \mathcal{S}_i(\gamma_i(u, \dots, u)), \quad (51)$$

where  $\kappa_n(u) = \theta_n(u)$  and, for  $1 \leq i \leq n-1$ ,

$$\kappa_i(u) = \tilde{\theta}_i \left( \prod_{j=i+1}^n C_{i,j}(\kappa_j(u)), \gamma_i(u, \dots, u) \right). \quad (52)$$

**Remark 6.** As indicated, we have assumed that the  $M_{i,j}$ 's are independent. However, the results can easily be extended to the case with dependencies by introducing  $C_i(\mathbf{z})$  as the joint PGF of  $M_{i,1}, \dots, M_{i,n}$ . It can be verified that the only difference would be that in (50) and (52), the product needs to be replaced by this joint PGF. For example, in (52) one would obtain

$$\kappa_i(u) = \tilde{\theta}_i(C_i(1, 1, \dots, 1, \kappa_{i+1}(u), \kappa_{i+2}(u), \dots, \kappa_n(u)), \gamma_i(u, \dots, u)).$$

In fact, by taking  $C_i(\mathbf{z}) = p_{i,0} + \sum_{j=1}^n p_{i,j} z_j$  we obtain a polling system with customer routing, as studied in [4, 18], where  $p_{i,j}$  denotes the probability that a customer finishing service at  $Q_i$  is routed to  $Q_j$ . For that particular choice of  $C_i(\mathbf{z})$ , our result agrees with the cycle time LST found in [4, Proposition 3.1].

## B 1-limited service

We consider the 1-limited service discipline. The goal of this appendix is to analyze the limiting value of  $\text{Cov}[L^e, I]$  as  $n$  grows large and to demonstrate why this is difficult. Furthermore, we intend to provide arguments that support the conjecture that the variance of the scaled cycle time is of order  $1/n$ .

Define  $g := \mathbb{P}[L^b > 0]$ . Note that the average number of customers that join  $Q_1$  during a cycle is given by  $\mathbb{E}[C]\Lambda/n$ . In order for the system to be stable, this should also equal the fraction of cycles that the server serves a customer at  $Q_1$ . Therefore,

$$g = \mathbb{P}[L^b > 0] = \frac{\Lambda}{n} \mathbb{E}[C] = \Lambda s / (1 - \rho).$$

Note that, since switch-over times during an intervisit time are independent of  $L^e$ , by additivity of the covariance function

$$\text{Cov}[L^e, I] = \sum_{i=2}^n \text{Cov}[L^e, V_i].$$

Furthermore, we can write

$$\mathbb{E}[L^e V_i] = \mathbb{E}[L^e] \beta - \mathbb{P}[Q_i \text{ is found empty}] \mathbb{E}[L^e \mid Q_i \text{ is found empty}] \beta.$$

Notice that  $\mathbb{P}[Q_i \text{ is found empty}] = 1 - g$ . Using  $\mathbb{E}[V_i] = g\beta$ , we conclude

$$\text{Cov}[L^e, I] = \beta(1 - g) \sum_{i=2}^n (\mathbb{E}[L^e] - \mathbb{E}[L^e \mid Q_i \text{ is found empty}]). \quad (53)$$

Hence, it remains to find  $\mathbb{E}[L^e]$  and  $\mathbb{E}[L^e \mid Q_i \text{ is found empty}]$ . We now turn to this task.

(i) Determination of  $\mathbb{E}[L^e]$ .

Write  $\mathcal{F}(z)$  for the PGF of the joint queue lengths at the start of a visit to  $Q_1$ . For 1-limited service it is well known (see for example [20]) that the following relation holds

$$\begin{aligned} \mathcal{F}(z_1, \dots, z_n) &= \left( \frac{1}{z_n} (\mathcal{F}(z_n, z_1, \dots, z_{n-1}) - \mathcal{F}(0, z_1, \dots, z_{n-1})) \mathcal{B} \left( \sum_{j=1}^n \frac{\Lambda}{n} (1 - z_j) \right) \right. \\ &\quad \left. + \mathcal{F}(0, z_1, \dots, z_{n-1}) \right) \mathcal{S} \left( \sum_{j=1}^n \frac{\Lambda}{n} (1 - z_j) \right). \end{aligned} \quad (54)$$

Differentiating once gives

$$l_i = \begin{cases} l_{i+1} + g\rho/n + s\Lambda/n = l_{i+1} + g/n, & i \neq n, \\ l_1 - g(1 - \rho/n) + s\Lambda/n = l_1 - (1 - 1/n)g, & i = n. \end{cases} \quad (55)$$

Summing gives

$$\sum_{i=1}^n l_i = nl_1 - (n-1)g/2. \quad (56)$$

Setting  $z_i = z$  for all  $i$  in Equation (54), we find

$$\mathcal{F}(z, \dots, z) = \frac{\mathcal{S}(\Lambda(1-z))(z - \mathcal{B}(\Lambda(1-z)))}{z - \mathcal{S}(\Lambda(1-z))\mathcal{B}(\Lambda(1-z))} \mathcal{F}(0, z, \dots, z). \quad (57)$$

We define  $\mathcal{G}_i(z) := \mathcal{F}(1, \dots, 1, z, 1, \dots, 1)$  and  $\mathcal{G}_i^0(z) := \mathcal{F}(0, 1, \dots, 1, z, 1, \dots, 1)$ , where  $z$  is the  $i$ 'th argument and  $\mathcal{G}_1^0 := 1 - g$ , and let  $l_i^0 := \frac{\partial}{\partial z_i} \mathcal{G}_i^0(z) \Big|_{z=(1, \dots, 1)}$ . Differentiating (57) then gives

$$\sum_{i=1}^n l_i = \frac{1}{1-g} \sum_{i=1}^n l_i^0 + \frac{g}{1-g} + \frac{\Lambda^2}{(1-g)(1-\rho)} \left( s^{(2)} - s^2 + \frac{1}{2}(g\beta^{(2)} - s^{(2)}) \right). \quad (58)$$

For a third equation, summing Equation (54) over all  $i$ , we find

$$\begin{aligned} \sum_{i=1}^n \mathcal{G}_i(z) &= \frac{\mathcal{S}\left(\frac{\Lambda}{n}(1-z)\right)(1-\mathcal{B}\left(\frac{\Lambda}{n}(1-z)\right))}{1-\mathcal{S}\left(\frac{\Lambda}{n}(1-z)\right)\mathcal{B}\left(\frac{\Lambda}{n}(1-z)\right)} \sum_{i=1}^n \mathcal{G}_i^0(z) \\ &\quad - \left(1-\frac{1}{z}\right) \frac{\mathcal{S}\left(\frac{\Lambda}{n}(1-z)\right)\mathcal{B}\left(\frac{\Lambda}{n}(1-z)\right)}{1-\mathcal{S}\left(\frac{\Lambda}{n}(1-z)\right)\mathcal{B}\left(\frac{\Lambda}{n}(1-z)\right)} (\mathcal{G}_1(z) - \mathcal{G}_1^0(z)). \end{aligned}$$

Differentiating and simplifying gives

$$\sum_{i=1}^n l_i = \frac{1}{\rho(1-g)+g} \left( nl_1 + \sum_{i=1}^n l_i^0 \right) - \frac{\Lambda^2(s^{(2)} - 2s^2 + g\beta^{(2)}) - 2g\rho + 2ng}{2(\rho(1-g)+g)}. \quad (59)$$

Finally, solving Equations (56), (58) and (59), we find (see also [20])

$$\begin{aligned} \sum_{i=1}^n l_i &= \frac{1}{1-g} \left( \frac{\Lambda^2 s^{(2)} - 2\Lambda^2 s^2 + (n+1)\Lambda s}{2(1-\rho)} + \frac{\Lambda^3 s \beta^{(2)}}{2(1-\rho)^2} \right), \\ \sum_{i=2}^n l_i^0 &= \frac{(n-1)\Lambda s}{2(1-\rho)} = (n-1) \frac{g}{2}, \end{aligned}$$

and

$$l_i = \frac{g(2-g)}{2(1-g)} + \frac{1}{n} \frac{\Lambda^2(s^{(2)} - s^2) + \Lambda^2 g(\beta s + \beta^{(2)})}{2(1-\rho)(1-g)} - \frac{i-1}{n} g.$$

Additionally, we find

$$l_i^0 = \frac{g(2-g)}{2} - \frac{i-1}{n} g(1-g).$$

Since on average  $g$  customers receive service during a server visit and  $g\rho/n$  customers join the queue during a server visit, it is easily seen that

$$\mathbb{E}[L^e] = l_1 - (1-\rho/n)g = \frac{g^2}{2(1-g)} + \frac{1}{n} \left( \frac{\Lambda^2(s^{(2)} - s^2) + \Lambda^2 g(\beta s + \beta^{(2)})}{2(1-\rho)(1-g)} + \rho g \right). \quad (60)$$

(ii) Determination of  $\mathbb{E}[L^e \mid Q_i \text{ is found empty}]$ .

Consider now a special cycle  $C^0 = \sum_{i=1}^n (V_i^0 + S_i^0)$  in which the server at the end of the cycle returns to find  $Q_1$  empty. Then, by symmetry, we can write

$$\mathbb{E}[L^e \mid Q_{n-i+2} \text{ is found empty}] = \frac{1}{1-g} l_i^0 - \frac{\Lambda}{n} \mathbb{E} \left[ S_i^0 + \sum_{j=i+1}^n (V_j^0 + S_j^0) \right].$$

Now, write

$$\mathbb{E}[S_i^0] = s - s^0(i) \text{ and } \mathbb{E}[V_i^0] = g\beta - v^0(i).$$

Then combining all results we find

$$\begin{aligned} \text{Cov}[L^e, I] &= \frac{n-1}{2n} \beta \left( \frac{\Lambda^2(s^{(2)} - s^2 + g\beta^{(2)})}{1-\rho} + \rho g^2 \right) \\ &\quad - \beta(1-g) \frac{\Lambda}{n} \left( \sum_{i=2}^n (i-1)s^0(i) + \sum_{j=3}^n (j-2)v^0(j) \right), \end{aligned}$$

and, taking limits and simplifying even further,

$$\begin{aligned} \lim_{n \rightarrow \infty} n\text{Var}[C/n] &= s^{(2)} - s^2 + g\beta^{(2)} - g^2\beta^2 \\ &+ (1-g) \lim_{n \rightarrow \infty} \frac{2}{n} \left( \sum_{i=2}^n (i-1)s^0(i) + \sum_{j=3}^n (j-2)v^0(j) \right). \end{aligned}$$

Disregarding the last term, this result already has the desired property that, as  $g \uparrow 1$ , the scaled variance converges to  $\text{Var}[S] + \text{Var}[B]$  and, as  $g \downarrow 0$ , the scaled variance converges to  $\text{Var}[S]$ .

It remains to determine  $s^0(i)$  and  $v^0(j)$ , which turns out to be non-trivial. However, since both  $s^0(i)$  and  $v^0(j)$  are most likely of order  $1/n$ , we conjecture that the variance of the scaled cycle time will also be of order  $1/n$ . To see this, we give the following heuristic. Ignoring any arrivals during  $S_i$  to queues other than  $Q_1$  whose service may cause an arrival to  $Q_1$ , we can write

$$E[S_i^0] \approx E[S_i | \text{no arrivals to } Q_1 \text{ during } S_i] = -\mathcal{S}'(\Lambda/n) = s - \Lambda s^{(2)}/n + O(1/n^2).$$

Hence,  $s^0(i) \approx \Lambda s^{(2)}/n$ . A similar argument can be applied to  $E[V_i^0]$ . As a consequence, also for 1-limited service, the scaled cycle time will converge in probability to a deterministic value, which we can exploit to determine the steady-state queue length distribution at visit beginnings.

## References

- [1] N. T. J. Bailey. On queueing processes with bulk service. *Journal of the Royal Statistical Society. Series B (Methodological)*, 16(1):80–87, 1954.
- [2] M. A. A. Boon, I. J. B. F. Adan, and O. J. Boxma. A polling model with multiple priority levels. *Performance Evaluation*, 67:468–484, 2010.
- [3] M. A. A. Boon, R. D. van der Mei, and E. M. M. Winands. Applications of polling systems. *Surveys in Operations Research and Management Science*, 16(2):67–82, 2011.
- [4] M. A. A. Boon, R. D. van der Mei, and E. M. M. Winands. Waiting times in queueing networks with a single shared server. *Queueing Systems*, 74(4):403–429, 2013.
- [5] S. C. Borst and O. J. Boxma. Polling: Past, present and perspective. *To appear in TOP*, 2018.
- [6] O. J. Boxma and W. P. Groenendijk. Pseudo-conservation laws in cyclic-service systems. *Journal of Applied Probability*, 24(4):949–964, 1987.
- [7] O. J. Boxma, J. Bruin, and B. Fralix. Sojourn times in polling systems with various service disciplines. *Performance Evaluation*, 66(11):621–639, 2009.
- [8] S. T. Bushby. BACnetTM: A standard communication infrastructure for intelligent buildings. *Automation in Construction*, 6(5):529–540, 1997.
- [9] G. L. Choudhury and W. Whitt. Computing distributions and moments in polling models by numerical transform inversion. *Performance Evaluation*, 25(4):267–292, 1996.
- [10] E. Coffman and E. Gilbert. A continuous polling system with constant service times. *IEEE Transactions on Information Theory*, 32(4):584–591, 1986.
- [11] S. W. Fuhrmann and R. B. Cooper. Stochastic decompositions in the M/G/1 queue with generalized vacations. *Operations Research*, 33(5):1117–1129, 1985.
- [12] N. K. Jaiswal. A bulk-service queueing problem with variable capacity. *Journal of the Royal Statistical Society. Series B (Methodological)*, 23(1):143–148, 1961.

- [13] W. Kastner, G. Neugschwandtner, S. Soucek, and H. Newmann. Communication systems for building automation and control. *Proceedings of the IEEE*, 93(6):1178–1203, 2005.
- [14] D. P. Kroese and V. Schmidt. A continuous polling system with general service times. *Annals of Applied Probability*, 2(4):906–927, 1992.
- [15] H. Levy. Binomial-gated service: a method for effective operation and optimization of polling systems. *IEEE Transactions on Communications*, 39(9):1341–1350, 1991.
- [16] H. Levy, M. Sidi, and O. J. Boxma. Dominance relations in polling systems. *Queueing Systems*, 6(1):155–171, 1990.
- [17] J. A. C. Resing. Polling systems and multitype branching processes. *Queueing Systems*, 13(4):409–426, 1993.
- [18] M. Sidi, H. Levy, and S. W. Fuhrmann. A queueing network with a single cyclically roving server. *Queueing Systems*, 11:121–144, 1992.
- [19] N. C. Strole. The IBM token-ring network - A functional overview. *IEEE Network*, 1(1):23–30, 1987.
- [20] H. Takagi. Mean message waiting times in symmetric multi-queue systems with cyclic service. *Performance Evaluation*, 5(4):271–277, 1985.
- [21] R. D. van der Mei and H. Levy. Polling systems in heavy traffic: Exhaustiveness of service policies. *Queueing Systems*, 27(3–4):227–250, 1997.
- [22] S. Wang. *Intelligent Buildings and Building Automation*. Taylor & Francis, New York, 2009.
- [23] E. M. M. Winands, I. J. B. F. Adan, and G.-J. van Houtum. Mean value analysis for polling systems. *Queueing Systems*, 54:35–44, 2006.